

**Учебно-исследовательская лаборатория
"Математические и программные технологии для
современных компьютерных систем
(Информационные технологии)"**

**Учебный курс
«Модели и методы
конечномерной оптимизации»**

**Часть 2
«Нелинейное
программирование и
многоэкстремальная
оптимизация»**

Оглавление

Глава 1. Математические модели оптимального выбора Ошибка! Закладка не определена.

1.1. Объект и его описание, модель процесса рационального выбора и постановки оптимизационных задач... **Ошибка! Закладка не определена.**

1.2. Различные трактовки решения в однокритериальных задачах, примеры **Ошибка! Закладка не определена.**

1.3. Понятия оптимальности в многокритериальных задачах и схемы компромисса **Ошибка! Закладка не определена.**

1.3.1. Концепции решений по Парето и Слейтеру .. **Ошибка! Закладка не определена.**

1.3.2. Лексикографическая схема компромисса **Ошибка! Закладка не определена.**

1.3.3. Метод главного критерия **Ошибка! Закладка не определена.**

1.3.4. Метод уступок **Ошибка! Закладка не определена.**

1.3.5. Метод идеальной точки Вержбицкого **Ошибка! Закладка не определена.**

1.3.6. Метод линейной свертки **Ошибка! Закладка не определена.**

1.3.7. Свертка Ю.Б. Гермейера..... **Ошибка! Закладка не определена.**

1.3.8 Проблема оценивания всего множества эффективных точек. . **Ошибка! Закладка не определена.**

1.4. Модели функций, используемые в задачах оптимального выбора **Ошибка! Закладка не определена.**

1.4.1. Модели функций, основанные на представлениях о выпуклости **Ошибка! Закладка не определена.**

1.4.1.1. Выпуклые, строго и сильно выпуклые функции **Ошибка! Закладка не определена.**

1.4.1.2. Квазивыпуклые, строго и сильно квазивыпуклые функции **Ошибка! Закладка не определена.**

1.4.1.3. Псевдовыпуклые и строго псевдовыпуклые функции **Ошибка! Закладка не определена.**

1.4.2. Модели функций используемые в многоэкстремальной оптимизации.... **Ошибка! Закладка не определена.**

1.4.2.1. Примеры детерминированных моделей многоэкстремальных функций **Ошибка! Закладка не определена.**

1.4.2.2. Примеры вероятностных моделей многоэкстремальных функций .. **Ошибка! Закладка не определена.**

1.4.2.3. Неполные адаптивные вероятностные модели **Ошибка! Закладка не определена.**

Глава 2. Теоретические основы аналитического решения задач оптимизации Ошибка! Закладка не определена.

2.1. Обобщение условий экстремума на задачи векторной оптимизации **Ошибка! Закладка не определена.**

2.2. Условия оптимальности в дифференциальной форме для многокритериальных задач оптимизации специального и общего вида **Ошибка! Закладка не определена.**

2.2.1. Условия первого порядка **Ошибка! Закладка не определена.**

2.2.2. Условия экстремума второго порядка **Ошибка! Закладка не определена.**

2.3. Элементы теории двойственности в задачах математического программирования с одним критерием **Ошибка! Закладка не определена.**

Глава 3. Общие методы учета ограничений в задачах математического программирования Ошибка! Закладка не определена.

- 3.1. Общие методы учета ограничений, обзор методов **Ошибка! Закладка не определена.**
- 3.2. Метод внешнего штрафа **Ошибка! Закладка не определена.**
 - 3.2.1. Общее описание и некоторые свойства **Ошибка! Закладка не определена.**
 - 3.2.2. Исследование сходимости и алгоритм настройки коэффициента штрафа **Ошибка! Закладка не определена.**
 - 3.2.3. Структура возникающих задач со штрафом и характер приближения оценок к решению **Ошибка! Закладка не определена.**
 - 3.2.4. Оценки скорости сходимости метода внешнего штрафа ... **Ошибка! Закладка не определена.**
 - 3.2.5. Недостаточность локальных методов при использовании метода штрафов **Ошибка! Закладка не определена.**
- 3.3. Метод модифицированных функций Лагранжа **Ошибка! Закладка не определена.**
 - 3.3.1. Общая схема метода множителей Лагранжа и ее недостатки. **Ошибка! Закладка не определена.**
 - 3.3.2. Преобразование постановки задачи, сведение задач с неравенствами к задачам с равенствами **Ошибка! Закладка не определена.**
 - 3.3.3. Построение модифицированной функции Лагранжа **Ошибка! Закладка не определена.**
 - 3.3.4. Метод модифицированной функции Лагранжа для задач с ограничениями– равенствами **Ошибка! Закладка не определена.**
 - 3.3.5. Метод модифицированной функции Лагранжа в задачах с равенствами и неравенствами **Ошибка! Закладка не определена.**
- 3.4. Другие общие методы учета ограничений **Ошибка! Закладка не определена.**
 - 3.4.1. Метод параметризации целевой функции **Ошибка! Закладка не определена.**
 - 3.4.2. Метод допустимой точки **Ошибка! Закладка не определена.**
 - 3.4.3. Индексный метод учета ограничений **Ошибка! Закладка не определена.**

Глава 4. Математические основы построения и анализа алгоритмов оптимизации Ошибка! Закладка не определена.

- 4.1. Модели численных методов оптимизации **Ошибка! Закладка не определена.**
 - 4.1.1. Основные обозначения **Ошибка! Закладка не определена.**
 - 4.1.2. Формальная модель и общая вычислительная схема **Ошибка! Закладка не определена.**
 - 4.1.3. Сходимость и оценки решения **Ошибка! Закладка не определена.**
- 4.2. Принципы построения методов оптимизации ... **Ошибка! Закладка не определена.**
- 4.3. Одношагово-оптимальные методы оптимизации .. **Ошибка! Закладка не определена.**
 - 4.3.1. Принцип одношаговой оптимальности **Ошибка! Закладка не определена.**
 - 4.3.2. Метод ломаных как одношагово-оптимальный алгоритм . **Ошибка! Закладка не определена.**
 - 4.3.3. Информационно-статистический метод Р.Г.Стронгина **Ошибка! Закладка не определена.**
 - 4.3.4. Одношагово-оптимальный байесовский метод Х.Кушнера.... **Ошибка! Закладка не определена.**
 - 4.3.5. Одношагово–оптимальный метод на основе адаптивных вероятностных моделей для задач с ограничениями **Ошибка! Закладка не определена.**
 - 4.3.6. Асимптотическая оптимальность **Ошибка! Закладка не определена.**
- 4.4. Теоретические основы сходимости одномерных алгоритмов глобального поиска **Ошибка! Закладка не определена.**

4.5. Анализ сходимости многомерных методов многоэкстремальной оптимизации. T–представимые алгоритмы и их свойства..... **Ошибка! Залка не определена.**

Залка не определена.

4.5.1. Описание класса задач **Ошибка! Залка не определена.**

4.5.2. Класс T–представимых алгоритмов, классификация, примеры..... **Ошибка!**

Залка не определена.

4.5.3. Теория сходимости T–представимых алгоритмов **Ошибка! Залка не определена.**

4.6. Анализ относительной плотности размещения испытаний при всюду плотной сходимости **Ошибка! Залка не определена.**

4.6.1. Необходимые понятия и обозначения..... **Ошибка! Залка не определена.**

4.6.2. Метод аналитического оценивания относительной концентрации испытаний **Ошибка! Залка не определена.**

Глава 5. Фундаментальные способы редукции размерности в многоэкстремальных задачах Ошибка! Залка не определена.

5.1. Многоэкстремальные задачи и методы покрытий. **Ошибка! Залка не определена.**

5.2. Принципы редукции сложности в многомерных многоэкстремальных задачах..... **Ошибка! Залка не определена.**

5.3. Многошаговая схема редукции размерности.... **Ошибка! Залка не определена.**

5.4. Свойства одномерных подзадач многошаговой схемы..... **Ошибка! Залка не определена.**

5.4.1. Структура допустимых областей одномерного поиска **Ошибка! Залка не определена.**

5.4.2. Свойства целевых функций в одномерных подзадачах..... **Ошибка! Залка не определена.**

5.5. Редукция размерности на основе кривых Пеано ... **Ошибка! Залка не определена.**

5.6. Решение задач с ограничениями с использованием разверток **Ошибка! Залка не определена.**

5.7. Компонентные методы **Ошибка! Залка не определена.**

5.7.1. Метод деления на три **Ошибка! Залка не определена.**

5.7.2. Диагональные компонентные методы Я.Пинтера..... **Ошибка! Залка не определена.**

5.7.3. Эффективные диагональные компонентные методы на основе адаптивных диагональных кривых **Ошибка! Залка не определена.**

5.7.4. Компонентные методы, основанные на триангуляции области поиска... **Ошибка!**

Залка не определена.

Глава 6. Методы построения оценок множества слабо эффективных точек, не использующие параметрических сверток .. Ошибка! Залка не определена.

6.1. Основные принципы непараметрической скаляризации **Ошибка! Залка не определена.**

6.1.1. Метод сведения к скалярной задаче с перестраиваемой целевой функцией **Ошибка! Залка не определена.**

6.1.2. Метод неравномерных покрытий Ю.Г. Евтушенко и М.А. Потапова **Ошибка! Залка не определена.**

6.1.3. Точное сведение многокритериальной задачи к скалярной с помощью свертки Д. Л. Маркина, Р. Г. Стронгина..... **Ошибка! Залка не определена.**

6.2. Реализация метода неравномерных покрытий Ю.Г. Евтушенко по схеме деления на три **Ошибка! Залка не определена.**

6.3. Одношагово–оптимальный метод многокритериальной оптимизации на основе адаптивных стохастических моделей.....**Ошибка! Закладка не определена.**

6.4. Метод построения равномерной оценки множества слабо эффективных точек..... **Ошибка! Закладка не определена.**

Глава 7. Модели и методы поиска локально–оптимальных решений
..... **Ошибка! Закладка не определена.**

7.1. Применение принципов оптимальности при построении методов локальной оптимизации выпуклых гладких задач...**Ошибка! Закладка не определена.**

7.1.1. Метод центров тяжести **Ошибка! Закладка не определена.**

7.1.2. Метод эллипсоидов **Ошибка! Закладка не определена.**

7.2. Принципы построения методов локальной оптимизации в задачах общего вида **Ошибка! Закладка не определена.**

7.2.1. Общая структура методов поиска локального минимума, принцип локального спуска **Ошибка! Закладка не определена.**

7.2.2. Измерения локальной информации и роль модели задачи в их интерпретации

..... **Ошибка! Закладка не определена.**

7.2.3. Классификация траекторных методов локального поиска. **Ошибка! Закладка не определена.**

7.2.4. Эффективные стратегии поиска вдоль направлений. Регуляризованные алгоритмы одномерного поиска **Ошибка! Закладка не определена.**

7.3. Аппроксимационные принципы построения алгоритмов. Анализ свойств классического градиентного метода и метода Ньютона ..**Ошибка! Закладка не определена.**

7.4. Эффективные методы второго порядка для гладких задач.....**Ошибка! Закладка не определена.**

7.4.1. Расширение области сходимости метода Ньютона за счет регулировки величины шага **Ошибка! Закладка не определена.**

7.4.2. Стратегии модификации матриц Гессе при нарушении их положительной определенности **Ошибка! Закладка не определена.**

7.5. Методы первого порядка, явно изменяющие метрику пространства
..... **Ошибка! Закладка не определена.**

7.5.1. Квазиньютоновские методы. Рекуррентные соотношения для оценок матриц Гессе по измерениям градиента в основных точках поиска..... **Ошибка! Закладка не определена.**

7.5.2. Модифицированные квазиньютоновские методы **Ошибка! Закладка не определена.**

7.5.3. Методы растяжения пространства (R–алгоритмы Н.З. Шора) **Ошибка! Закладка не определена.**

7.6. Методы сопряженных направлений**Ошибка! Закладка не определена.**

7.6.1. Сопряженные направления и их свойства **Ошибка! Закладка не определена.**

7.6.2. Метод сопряженных градиентов Флетчера-Ривса **Ошибка! Закладка не определена.**

7.7. Некоторые методы прямого поиска для негладких задач**Ошибка! Закладка не определена.**

7.7.1. Метод Нелдера–Мида..... **Ошибка! Закладка не определена.**

7.7.2. Метод Хука-Дживса..... **Ошибка! Закладка не определена.**

7.8. Специальные методы учета линейных ограничений в гладких задачах локальной оптимизации **Ошибка! Закладка не определена.**

7.8.1. Специальные методы учета линейных равенств **Ошибка! Закладка не определена.**

7.8.2. Специальные методы учета линейных неравенств, методы активного набора **Ошибка! Закладка не определена.**

7.8.3. Особенности применения методов локального поиска при двусторонних ограничениях на переменные..... **Ошибка! Закладка не определена.**

 7.8.3.1 Особенности учета двусторонних ограничений на переменные в методах гладкой оптимизации **Ошибка! Закладка не определена.**

 7.8.3.2. Учет двусторонних ограничений в методах прямого поиска..... **Ошибка! Закладка не определена.**

Заключение **Ошибка! Закладка не определена.**

Литература **Ошибка! Закладка не определена.**

 Основная литература **Ошибка! Закладка не определена.**

 Дополнительная литература **Ошибка! Закладка не определена.**

[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_1.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_2.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_3.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_4.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_5.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_6.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optChapter_2_7.doc](#)
[D:\OPT_Optimisation\02_np&mo\Book\optReferences_2.doc](#)

Глава 1. Математические модели оптимального выбора

1.1. Объект и его описание, модель процесса рационального выбора и постановки оптимизационных задач

Книгу по теории и вычислительным методам конечномерной не дискретной оптимизации стоит начать с изложения взгляда авторов на то место, которое занимают подобные задачи по отношению к более общей проблеме рационального выбора¹.

При решении реальных прикладных задач, связанных с рациональным выбором, обычно приходится сталкиваться с отсутствием их формальной математической постановки. Как правило, имеется лишь некоторый объект по отношению к параметрам P которого нужно сделать выбор, а также набор числовых характеристик–критериев $W(P)$, определяющих зависимость количественных показатели свойства объекта от значений его параметров. Эти критерии можно назвать характеристиками функционирования объекта. На основе их значений должен происходить выбор параметров P в пределах некоторых ограничений $A \leq P \leq B$, вытекающих из смысла решаемой задачи. При этом объект обычно представлен своей вычислительной моделью (рис. 1.1).



Рис. 1.1. Основные компоненты математической модели объекта

Постановщик задачи обычно обладает некоторой, точно не формализуемой информацией о целях рационального выбора. Причем цели могут изменяться в зависимости от дополнительной информации, получаемой постановщиком в процессе решения, а это, в свою очередь, может приводить к изменению используемой модели объекта (рис.1.2).

На каждом этапе решения интуитивно понимаемая постановщиком цель формализуется, с участием второго лица — математика–вычислителя, в виде вспомогательной задачи оптимального выбора. Результат ее решения может приводить к изменению постановки этой задачи или даже к изменению принятой модели объекта.

¹ В этом разделе авторы используют переработанный материал из их книги [43].

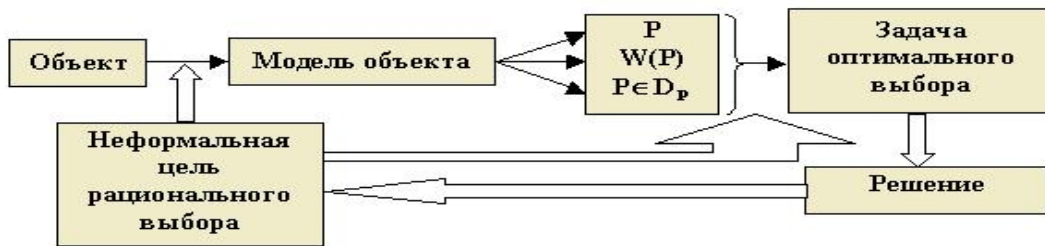


Рис. 1. 2. Соотношение неформальной цели и задач оптимального выбора

Конкретная постановка задачи оптимального выбора связана с фиксацией части параметров, установкой ограничителей на переменные, оставшиеся свободными (рис.1.3),



Рис. 1. 3. Разделение исходных параметров объекта на переменные и константы

выделением групп оптимизируемых и ограничиваемых характеристик функционирования (рис. 1.4).



Рис. 1. 4. Перераспределение критериев по группам

Множества J_f , J_g , J_h содержат номера критериев функционирования, соответствующих компонентам вектор-функций f , g , h . Кроме того, на значения функций g и h накладываются свои ограничители g^+ и \bar{h} : $g(y) \leq g^+$, $h(y) = \bar{h}$. Чтобы не усложнять запись в дальнейшем, мы будем переносить ограничители в левую часть неравенств и равенств, пряча их в функциях g и h . В результате ограничения примут вид $g(y) \leq 0$, $h(y) = 0$.

Таким образом, выбор $\langle J_c, J_y, C, a, b, J_f, J_g, J_h, g^+, \bar{h} \rangle$ породит экстремальную задачу оптимального выбора.

$$f(y) \rightarrow \inf, y \in Y, f : D \rightarrow R^n, \quad (1.1)$$

$$Y = \{y \in D \subseteq R^N; g(y) \leq 0, h(y) = 0\}, \quad (1.2)$$

$$g : D \rightarrow R^m, h : D \rightarrow R^p,$$

$$D = \{y \in R^N : a \leq y \leq b\}, \quad (1.3)$$

где f – целевая вектор-функция, а Y – допустимая область.

Задачи такого вида могут быть решены с использованием теории и вычислительных методов конечномерной оптимизации. Именно этим вопросам будем посвящен последующий материал данной книги. Следует подчеркнуть, что решение отдельной экстремальной задачи, вообще говоря, не является решением общей неформальной задачи рационального выбора, а лишь этапом ее решения. Рациональный выбор параметров объекта обычно приводит к решению последовательности экстремальных задач (1.1)-(1.3), вид и количество которых определяется поставщиком общей задачи совместно с математиком-вычислителем.

Процесс формулировки экстремальных задач является нетривиальным и не формализуемым элементом всего процесса решения и не является математической задачей, хотя и требует от сторон, принимающих решение, достаточной квалификации в области экстремальных задач. При постановке экстремальных задач (1.1)-(1.3) всегда следует обращать самое пристальное внимание на свойства полученных задач, поскольку они определяют в дальнейшем возможность их решения вычислительными методами или же аналитического решения.

В дальнейшем всегда будем предполагать, что задача оптимального выбора (1.1)-(1.3) уже поставлена. Нам предстоит формализовать понятие решения экстремальных задач (это, как будет показано, можно сделать не единственным образом). Далее будет построена теория, предоставляющая инструменты решения поставленных задач. Везде, где это возможно, изложение будет доведено до уровня построения аналитических либо вычислительных методов поиска решения.

1.2. Различные трактовки решения в однокритериальных задачах, примеры

Прежде чем переходить к определению понятия решения в многокритериальных задачах (1.1)–(1.3) с $n > 1$, рассмотрим важный частный случай $n=1$, т.е. задачи с одним целевым критерием (скалярные задачи). Их обычно называют *задачами математического программирования*. Введем ряд обозначений и напомним терминологию.

Определение 1.1. Если существует точка $y^* \in Y$, что

$$f(y^*) = \inf \{f(y) : y \in Y\}, \quad (1.4)$$

то говорят, что f достигает своей точной нижней грани на Y , а y^* называют точкой глобального (абсолютного) минимума, значение $f^* = f(y^*)$ называют наименьшим или глобально-оптимальным (минимальным) значением f на Y .

В некоторых случаях будет использоваться запись,

$$y^* = \arg \min \{f(y) : y \in Y\} \quad (1.5)$$

Множество всех точек $y^* \in Y$, удовлетворяющих (1.4) будем обозначать через

$$Y^* = \text{Arg} \min \{f(y) : y \in Y\} = \{y^* \in Y : f(y) = f^*\} \quad (1.6)$$

Определение 1.2. Точка $y^o \in Y$ называется точкой локального минимума f на Y , если существует такое $\varepsilon > 0$, что $\forall y \in Y \cap O_\varepsilon(y^o)$, где $O_\varepsilon(y^o) = \{y \in R^N : \|y - y^o\| < \varepsilon\}$, выполняется неравенство $f(y^o) \leq f(y)$.

Даже в случае однокритериальных задач могут использоваться разные трактовки понятия «решение». Выбор зависит, как правило, не от прихоти

постановщика, а от существования решаемой задачи. Ниже различные понятия решения сопровождаются примерами адекватных им задач.

Постановка А. Решение задачи (1.1)-(1.3) при $n=1$ заключается в определении экстремального значения

$$f^* = \inf \{f(y) : y \in Y\}. \quad (1.7)$$


Постановка В. Определить нижнюю грань из (1.7) и, если множество точек глобального минимума из (1.6) не пусто, найти хотя бы одну точку $y^* \in Y^*$.

Постановка С. Найти нижнюю грань f^* из (1.7) и определить все точки глобального минимума (либо убедиться, что множество Y^* пусто).

Постановка D. Найти точки y^o и значения $f(y^o)$ всех локальных минимумов.

Постановка Е. Для $y \in Y$ найти точку $y^o(y)$ и значение того локального минимума, в области притяжения которого лежит заданная начальная точка y . Последнее требование означает, что существует непрерывная кривая, проходящая в области Y от y к $y^o(y)$ при движении вдоль которой значение функции f не возрастает.

Постановки А–Е задачи математического программирования являются наиболее распространенными, хотя возможны и другие варианты. Искомые экстремальные характеристики, определяемые постановками А, В, ..., Е назовем А, В, ... или Е–*решениями задачи математического программирования* (1.1)-(1.3) ($n=1$).

 **Замечание.** Задача математического программирования (1.1)-(1.3) при $n=1$, для которой заведомо известно, что множество точек глобального минимума не пусто, будем записывать как

$$f(y) \rightarrow \min, y \in Y. \quad (1.8)$$

Достаточные условия не пустоты Y даются, например, теоремой Вейерштрасса.

Рассмотрим примеры задач, соответствующих основным постановкам А–Е.

Пример А. Задача вычисления меры робастной устойчивости по линейно входящим параметрам.

Пусть имеется динамическая система с непрерывным временем, имеющая характеристическое уравнение вида $\chi(p, a) = a_s \varphi_s(p) + \dots + a_1 \varphi_1(p) = 0$, где a – параметры системы, а p – комплексная переменная. Пусть U – область устойчивости системы в пространстве переменных a (при $a \in \text{int } U$ все корни p_k характеристического уравнения имеют $\text{Re } p_k < 0$). Пусть, кроме того, задана точка a^* – параметры конкретной динамической системы из рассматриваемого семейства, причем, $a^* \in U$. Требуется определить ρ^* -меру робастной устойчивости системы с параметрами a^* , т.е.

$$\rho^* = \inf \{ \|a - a^*\| : a \in \partial U \}, \quad (1.9)$$

где ∂U – границы U .

Поскольку уравнение границы может быть записано в виде параметрической системы

$$\begin{aligned} h_1(a, \omega) &= a_s \varphi_{s1}(\omega) + \dots + a_1 \varphi_{11}(\omega) = 0 \\ h_2(a, \omega) &= a_s \varphi_{s2}(\omega) + \dots + a_1 \varphi_{12}(\omega) = 0, \end{aligned} \quad (1.10)$$

$\omega \in (-\infty, +\infty)$, где $\varphi_k(i\omega) = \varphi_{k1}(\omega) + i\varphi_{k2}(\omega)$ ($k = 1, \dots, s$), то получаем для евклидовой метрики экстремальную задачу


$$f(a, \omega) = \sum_{i=1}^s (a_i - a_i^*)^2 \rightarrow \inf, \quad (a, \omega) \in Y, \quad (1.11)$$

$$Y = \{ (a, \omega) \in R^{s+1} : \omega \in R^1, h_i(a, \omega) = 0, (i=1,2) \}. \quad (1.12)$$

В этой задаче требуется найти только значение

$$\rho^* = (\inf\{f(a, \omega) : (a, \omega) \in Y\})^{0.5},$$

другие экстремальные характеристики не нужны, т.е. она соответствует постановке А — нужно найти А-решение.

 **Замечание.** Квадратичность по a целевой функции и линейность по a ограничений позволяют при фиксированном ω аналитически найти верхнюю грань по a в (1.11), используя условия экстремума. При этом будут получены значения $a_i(\omega)$ и задача сведется к виду $f(a(\omega), \omega) \rightarrow \inf, \omega \in (-\infty, +\infty)$, где целевая функция обычно является многоэкстремальной по ω .

Примеры В, С. Задача определения обобщенных координат манипулятора, обеспечивающих требуемое положение и ориентацию схвата.

Допустим, разрабатывается математическое обеспечение для управления рукой робота-манипулятора с шестью степенями свободы. Пусть известна вектор-функция $F(y)$, определяющая решение прямой кинематической задачи для схвата манипулятора. $y = (y_1, \dots, y_6)$ – вектор обобщенных координат манипулятора, первые три компоненты $(F_1(y), F_2(y), F_3(y))$ – определяют координаты характерной точки схвата, а последние три $(F_4(y), F_5(y), F_6(y))$ – три угла его ориентации.

Обратная кинематическая задача требует определения значений обобщенных координат y^* из области возможных значений $y^* \in D$, обеспечивающих заданное положение и ориентацию схвата, определяемые вектором F^* . Задача сводится к экстремальной

$$f(y) = \|F(y) - F^*\| \rightarrow \min, y \in Y = D \subset R^N, \quad N = 6.$$

Если дополнительно задана функция $\rho(y)$ — расстояние манипулятора до препятствий, то допустимая область примет вид $Y = \{y \in D : \rho(y) > \delta\}$, где $\delta > 0$ – область безопасности. Данная задача соответствует постановке В или С, поскольку требует определения координат точки (точек) глобального минимума. Если при этом окажется что $f(y^*) \neq 0$, то обратная кинематическая задача не имеет решений, т.е. схват манипулятора не может быть перемещен в точку (F_1^*, F_2^*, F_3^*) с ориентацией (F_4^*, F_5^*, F_6^*) . Знание всех точек глобального минимума может потребоваться в этой задаче для того, чтобы выбрать наилучшую конфигурацию манипулятора из возможных $y^* \in Y$ по какому-либо дополнительному критерию, например, по наибольшей близости к заданной текущей конфигурации. В этом случае задача соответствует постановке С.

Пример D. Задача местоопределения по измерениям рельефа местности.

Данная задача относится к группе задач навигации по геофизическим полям².

Рассмотрим математическую постановку задачи. Пусть имеется подвижный объект, равномерно и прямолинейно перемещающийся на постоянной высоте над участком поверхности Земли. С этим участком связана система координат y_1, y_2 . Закон изменения координат объекта следующий

$$y_1(t) = y_1 + v_1 t, \quad y_2(t) = y_2 + v_2 t. \quad (1.19)$$

Начальное местоположение (y_1, y_2) точно неизвестно и подлежит определению. Известной считается только область D его возможных значений

$$D = \{(y_1, y_2): a_1 \leq y_1 \leq b_1, a_2 \leq y_2 \leq b_2\}.$$

В известные моменты времени t_1, t_2, \dots, t_n подвижный объект определяет высоту рельефа местности в точках своего текущего местоположения, получая результаты измерений

$$h_i = h(y_1(t_i), y_2(t_i)) + C + \xi_i \quad (i=1, \dots, k) \quad (1.20)$$

где ξ_i – независимые реализации центрированной составляющей помехи измерений с плотностью распределения $P_\xi(z)$, C – систематическая составляющая помехи измерений, $h(y_1, y_2)$ – функция высоты рельефа местности. На борту объекта имеется электронная карта, позволяющая вычислять значения функции $h(y_1, y_2)$. Требуется по известной функции $h(y_1, y_2)$, значениям v_1, v_2 проекций вектора скорости, моментам времени t_1, t_2, \dots, t_k и результатам h_1, h_2, \dots, h_k измерений высот рельефа вдоль траектории полета оценить координаты y_1, y_2 начального местоположения объекта в области D .

Эта задача сводится к задаче оптимизации с использованием метода максимального правдоподобия. Будем считать плотность распределения помехи наблюдений известной. Тогда можно вычислить функцию $F(y_1, y_2, C)$, определяющую плотность вероятности наблюдаемых значений высот рельефа при условии, что начальное местоположение и систематическая составляющая измерений имеют значения y_1, y_2, C .

$$F(y_1, y_2, C) = P(h_1, h_2, \dots, h_k / y_1, y_2, C) = \prod_{i=1}^k P_\xi(h_i - h(y_1 + v_1 t_i, y_2 + v_2 t_i) - C). \quad (1.21)$$

Метод максимального правдоподобия сводится к определению оценок y_1, y_2 из решения экстремальной задачи

$$F(y_1, y_2, C) \rightarrow \max, \quad (y_1, y_2) \in D, \quad C \in R^1 \quad (1.22)$$

Наиболее простую форму в методе максимального правдоподобия задача приобретает в том случае, когда распределение ξ_i нормально. В этом случае

$$P_\xi(z) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{z^2}{2\sigma^2}}.$$

Легко видеть, что задача (1.22) преобразуется в форму задачи метода наименьших квадратов

² Подробнее об этих задачах смотрите в книге Белоглазов И.Н., Джанджагава Г.И., Чигин Г.П. Основы навигации по геофизическим полям /Под ред.А.А.Красовского. — М.: Наука, 1985.

$$f(y_1, y_2, C) \rightarrow \min, (y_1, y_2) \in D, -\infty < C < +\infty, \quad (1.23)$$

где

$$f(y_1, y_2, C) = 1/k \sum_{i=1}^k (h_i - h(y_1(t_i), y_2(t_i)) - C)^2 \quad (1.24)$$

При этом минимум по C можно найти аналитически из условия

$$\frac{\partial f}{\partial C}(y_1, y_2, C) = 0.$$

Отсюда находим

$$C = 1/k \sum_{j=1}^k (h_j - h(y_1(t_j), y_2(t_j))).$$

Окончательно, задача определения местоположения сводится к следующей задаче многоэкстремальной оптимизации

$$f(y_1, y_2) \rightarrow \min, (y_1, y_2) \in D, \quad (1.25)$$

$$f(y_1, y_2) = 1/k \sum_{i=1}^k \{ h_i - h(y_1 + v_1 t_i, y_2 + v_2 t_i) - 1/k \sum_{j=1}^k (h_j - h(y_1 + v_1 t_j, y_2 + v_2 t_j)) \}^2 \quad (1.26)$$

Многоэкстремальный характер функции (1.19) связан с существованием на карте местности участков с похожими сечениями рельефа. Решение при достаточно большом количестве измерений определяется глобальным минимумом. Если число измерений k невелико, то могут существовать локальные минимумы со значениями, близкими к глобальному, и тогда истинному местоположению объекта может соответствовать один из них. Для принятия обоснованного решения о местоположении, в этом случае нужна полная информация о положении и значениях всех локальных минимумов, т.е. в этой задаче необходимо определять D -решение.

Пример Е. Слежение за дрейфом минимума.

При решении задач математического программирования в присутствии ограничений одним из распространенных подходов к их учету является метод штрафов. Он будет подробно изучаться в главе 3. Сейчас мы используем вспомогательные задачи метода штрафов, как пример задач, решаемых в постановке Е. Действительно, в методе штрафов конструируются вспомогательные задачи

$$S_\gamma(y) = f(y) + \gamma H(y) \rightarrow \min, y \in D, \quad (1.20)$$

где $H(y)$ -функция штрафа, принимающая в Y значения равные 0, а вне Y — положительные значения. $H(y)$ вычисляется по функциям ограничений $g(y)$ и $h(y)$.

Должна быть решена последовательность таких задач для возрастающей последовательности γ_k , стремящейся к $+\infty$ при $k \rightarrow \infty$. Для получения глобального решения $y_{\gamma_{k+1}}^*$ ($k+1$)-й задачи обычно определяют ее локальный минимум, выбирая в качестве начальной точки поиска результата решения предыдущей задачи, т.е. определяют $y_{\gamma_{k+1}}^o(y_{\gamma_k}^*)$, находя E -решение задачи (1.20).

1.3. Понятия оптимальности в многокритериальных задачах и схемы компромисса

1.3.1. Концепции решений по Парето и Слейтеру

В этом пункте мы вернемся к общей многокритериальной постановке экстремальной задачи с ограничениями (1.1)–(1.3). Далее всегда будем предполагать, что каждая компонента $f_i(y)$ ($i=1, \dots, n$) достигает на допустимом множестве Y своей точной нижней грани. Перепишем задачу в следующем виде

$$f(y) = (f_1(y), \dots, f_n(y)) \rightarrow \min, \quad y \in Y, \quad (1.21)$$

$$Y = \{y \in D \subseteq R^N; g(y) \leq 0, h(y) = 0\}, \quad (1.22)$$

$$g(y) = (g_1(y), \dots, g_m(y)), \quad h(y) = (h_1(y), \dots, h_p(y)),$$

$$D = \{y \in R^N : a \leq y \leq b\}, \quad (1.23)$$

Вектор-функция f , определенная на D , переводит область $Y \subseteq R^N$ в некоторое множество $F = f(Y) \subseteq R^n$. Будем говорить, что Y –множество возможных (допустимых) значений переменных, а F –множество возможных оценок. Элементы из Y , как и раньше, будем обозначать через y , а элементы из F – через z (рис.1.5).

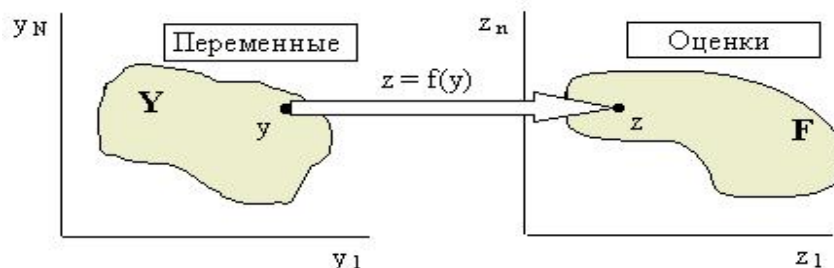
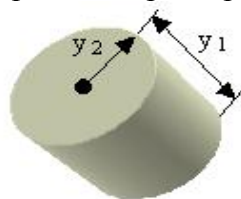


Рис. 1. 5. Множество допустимых значений переменных и множество возможных оценок

Пример. Рассмотрим простую задачу о максимизации объема цилиндрической емкости при минимизации площади ее поверхности с ограничением отношения геометрических размеров.



Инвертируя знак у критерия объема, и требуя выполнения нужного соотношения размеров, например, $y_2 \geq y_1/5$, получим задачу вида (1.21)–(1.23) с

$$f_1(y) = -\pi y_1 y_2^2,$$

$$f_2(y) = 2\pi y_2(y_2 + y_1),$$

$$g_1(y) = (y_1/5) - y_2, \quad n = 2, \quad N = 2, \quad m = 1$$

$$D = \{y = (y_1, y_2) : 0 \leq y_1 \leq 1, \quad 0 \leq y_2 \leq 1\}$$

На рис.1.6 представлены изолинии критериев, а также множество эффективных точек Y^* (определение дано ниже). Точки $y^{*1} \neq y^{*2}$ показывают

положение глобальных минимумов в каждой из компонент векторного критерия. Очевидна противоречивость требования одновременного достижения двух этих экстремальных значений.

Замечание. Общая ситуация характеризуется противоречивостью требования одновременной минимизации компонент векторного критерия. Многокритериальные задачи требуют новой концепции понятия решения.

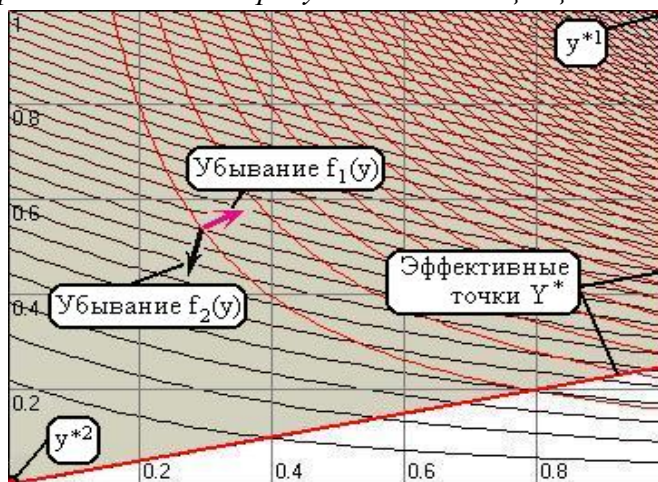


Рис. 1. 6. Поведение критериев и вид множества Y^* для примера с цилиндрической емкостью

Определение 1.3. Введем следующие отношения предпочтения. Будем говорить, что $f(y^1) \leq f(y^2)$, т.е. векторная оценка $f(y^1)$ не хуже $f(y^2)$, если $\forall i=1, \dots, n$ $f_i(y^1) \leq f_i(y^2)$. Будем говорить, что $f(y^1) < f(y^2)$, т.е. векторная оценка $f(y^1)$ лучше $f(y^2)$, если $\forall i=1, \dots, n$ $f_i(y^1) < f_i(y^2)$.

Заметим, что если следовать данному определению, то понятие «не хуже» не будет являться отрицанием понятия «хуже», трактуемым как покомпонентное выполнение неравенства $f(y^1) > f(y^2)$.

Введем в рассмотрение конус доминирования $R^-(z^o)$ для оценки z^o :

$$R^-(z^o) = \{z : z < z^o\} \subset R^n,$$

все точки которого лучше z^o , а также его замыкание $\bar{R}^-(z^o) = \{z : z \leq z^o\} \subset R^n$. В конусе (замкнутом конусе) доминирования, построенном для оценки z^o , находятся векторы z лучшие (не худшие) оценки z^o . Теперь можно естественным образом обобщить понятие оптимального решения на многокритериальные задачи.

Естественно считать оптимальными те векторы z^* из F , для которых в замыкании их конуса доминирования, т.е. в $\bar{R}^-(z^*)$ нет ни одного вектора из F , кроме самого z^* . Говорят, что оценка z^* оптимальна по Парето (эффективна). Для нее не найдется ни одной возможной оценки z из F , являющейся не худшей по сравнению с z^* (по каждой из компонент). Если смягчить условия и требовать лишь, чтобы не существовало лучших, чем z^o оценок в F , то оценку z^o называют оптимальной по Слейтеру или слабо эффективной.

Итак, мы приходим к следующим определениям оптимальности по Парето и Слейтеру.

Определение 1.4. Оценка z^* из F называется оптимальной по Парето (эффективной), если $\bar{R}(z^*) \cap F \setminus \{z^*\} = \emptyset$.

Оценка z^0 из F называется оптимальной по Слейтеру (слабо эффективной) оценкой, если $\bar{R}(z^0) \cap F = \emptyset$.

Точка y^* из Y называется эффективной (оптимальной по Парето), если ее оценка $f(y^*)$ оптимальна по Парето, т.е. не существует $y \in Y$ с $f(y) \neq f(y^*)$, что $f(y) \preceq f(y^*)$.

Точка y^0 из Y называется слабо эффективной (оптимальной по Слейтеру), если оценка $f(y^0)$ оптимальна по Слейтеру, т.е. не существует $y \in Y$, что $f(y) < f(y^0)$.

Как видно из определения, возможна определенная вариабельность в терминологии. В дальнейшем термин «эффективность» и «слабая эффективность» будем чаще относить к пространству параметров.

Обозначим через Y^* и Y^0 множества эффективных и слабо эффективных точек, а через $P(\cdot)$ и $S(\cdot)$ — операторы, выделяющие из множества F подмножества точек Парето $P(F)$ и Слейтера $S(F)$. Заметим, что всегда $P(F) \subseteq S(F)$.

Замечание. Следует обратить внимание на то, что в многокритериальных задачах оба понятия (как эффективных, так и слабо эффективных точек) обобщают понятие глобального минимума, используемое для однокритериальных задач.

Можно определить понятия локально эффективного решения и локально слабо эффективного решения.

Определение 1.5. Точку y^{*0} из Y назовем локально эффективной, если $\exists \varepsilon > 0$, что не существует точек $y \in O_\varepsilon(y^{*0}) \cap Y$ с $f(y) \neq f(y^{*0})$, что $f(y) \leq f(y^{*0})$.

Точку y^{00} из Y назовем локально слабо эффективной, если $\exists \varepsilon > 0$, что не существует $y \in O_\varepsilon(y^{00}) \cap Y$, что $f(y) < f(y^{00})$.

На рис.1.7, 1.8 приведены примеры, иллюстрирующие топологическую структуру множеств Y^* и Y^0 , а также $P(F)$ и $S(F)$ в конкретных ситуациях.

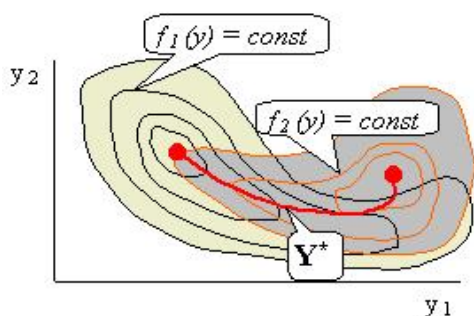


Рис.1.7. Случай $n=2, N=2$. Y^* — кривая

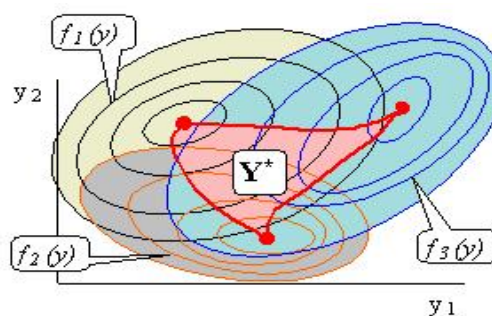


Рис.1.8. Случай $n=3, N=2$. Y^* — криволинейный треугольник

На рис.1.9 в пространстве критериев показана структура множества F для задачи с $f_1(y)=y$, $f_2(y)=\sin(y)$, $y \in [0, 4\pi]$. Пунктирной линией отмечены локально Парето-оптимальные оценки, а толстой линией — оптимальные по Парето.

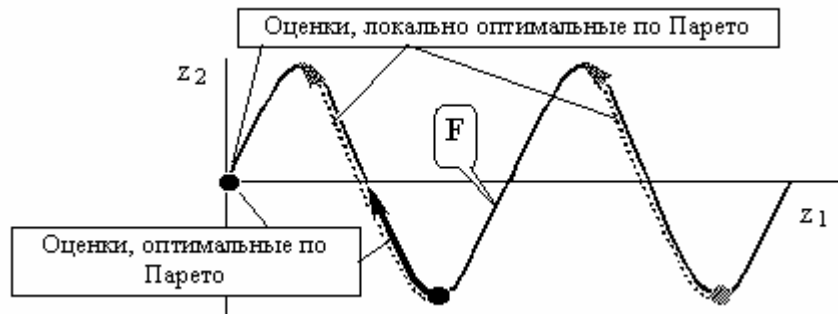


Рис.1. 9. Иллюстрация в пространстве критериев, случай $n=2, N=1$

Перейдем к способам оценивания решений многокритериальных задач. Возможно два различных подхода. Первый основан на использовании различных схем компромисса для сведения многокритериальной задачи ко множеству однокритериальных задач. Их решения соответствуют отдельным решениям исходной многокритериальной задачи. При использовании второго подхода строится оценка сразу всего множества Парето или Слейтера без сведения к отдельным скалярным задачам. Этот подход будет изложен в главе 6.

В этом разделе рассмотрим несколько схем компромисса, позволяющих выполнить переход к семейству однокритериальных задач.

1.3.2. Лексикографическая схема компромисса

Эта схема применяется в том случае, когда критерии упорядочены по важности:

f_1 – неизмеримо важнее f_2 , f_2 – неизмеримо важнее f_3 и т.д.

В этом случае задаче соответствует следующая схема решения. Обозначим $Y_0^* = Y$, определим

$$Y_1^* = \{y^{*1} \in Y_0^* : f_1(y^{*1}) = \min \{f_1(y) : y \in Y_0^*\}\},$$

$$Y_k^* = \{y^{*k} \in Y_{k-1}^* : f_k(y^{*k}) = \min \{f_k(y) : y \in Y_{k-1}^*\}\}$$

$(k = 2, \dots, n).$

Для случая двух критериев ($n=2$) процесс решения иллюстрируется на рис.1.10.

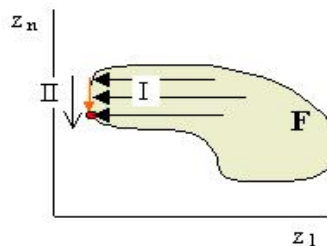


Рис. 1. 10. Поиск решения по лексикографической схеме

Свойство 1.1. Решения Y_n^* , найденные при лексикографической схеме компромисса принадлежат множеству эффективных решений, т.е. $Y_n^* \subseteq Y^*$.

ДОКАЗАТЕЛЬСТВО. Пусть это не так и $\exists \bar{y} \in Y, y^{*n} \in Y_n^*$, что $f(\bar{y}) \neq f(y^{*n})$ и $f(\bar{y}) \leq f(y^{*n})$. Пусть j – наименьший из индексов, при котором $f_j(\bar{y}) < f_j(y^{*n})$. Тогда $f_i(\bar{y}) = f_i(y^{*n}), (i=0, \dots, j-1)$, т.к. $y^{*n} \in Y_i^*$, т.е., $\bar{y} \in Y_{j-1}^*$. Но $Y_n^* \subseteq Y_j^*$, поэтому

для $y^{*n} \in Y_n^*$ будут выполняться неравенства $f_j(y^{*n}) \leq f_j(y) \quad \forall y \in Y_{j-1}^*$. Т.к. $\bar{y} \in Y_{j-1}^*$, то для него также выполнится неравенство $f_j(y^{*n}) \leq f_j(\bar{y})$.

Мы приходим к противоречию, поскольку для $f_j(\bar{y})$ должны одновременно выполняться два противоречивых неравенства.

1.3.3. Метод главного критерия

Метод заключается в том, что один из критериев объявляется главным, например, f_j . Он выбирается в качестве минимизируемого, а остальные критерии переводятся в категорию ограничений с верхними ограничителями f_i^+ . Возникает однокритериальная задача следующего вида

$$f_j(y) \rightarrow \min, y \in Y, f_i(y) \leq f_i^+ (i \neq j, i = 1, 2, \dots, n) \quad (1.24)$$

Свойство 1.2.

А. Если задача (1.24) имеет решение \bar{y} , то $\bar{y} \in Y^0$, т.е. является слабо эффективным.

В. Если y^ эффективная точка, то при $f_i^+ = f_i(y^*)$ y^* будет являться единственным (с точностью до f -эквивалентности) решением задачи (1.24).*

ДОКАЗАТЕЛЬСТВО. Обоснуем свойство (А). Предположим, что оно не верно, тогда найдется $y' \in Y$, что $f(y') < f(\bar{y})$. Если бы это было так, то все дополнительные неравенства были бы выполнены в точке y' , т.к. $f_i(y') < f_i(\bar{y}) \leq f_i^+ \quad \forall i \neq j$, и, следовательно, неравенство $f_j(y') < f_j(\bar{y})$ противоречило бы глобальной оптимальности \bar{y} в задаче (1.24).

Обоснуем (В). Поясним терминологию. f -эквивалентными называют две точки y', y'' из Y , если $f(y') = f(y'')$.

Поскольку $y^* \in Y^*$, то не найдется такой $y \in Y$, что $f(y) < f(y^*)$ и $f(y) \leq f(y^*)$. Поэтому на множестве

$$Y_j = \{y \in Y: f_i(y) \leq f_i(y^*) \quad i=1, \dots, n; i \neq j\}$$

критерий $f_j(y)$ будет достигать минимума в точке y^* и не будет других точек $y \in Y_j$ с $f(y) < f(y^*)$ что $f_j(y) \leq f_j(y^*)$. Таким образом, свойство (В) верно.

Следствие. При произвольном выборе главного критерия f_j всегда можно подобрать параметры f_i^+ для $i \neq j$, что решением задачи (1.24) будет любая эффективная точка.

Заметим, что основная проблема использования метода главного критерия состоит в том, что при произвольно выбранных ограничителях допустимая область в (1.24) может оказаться пустой.

1.3.4. Метод уступок

Этот подход по своей сути совпадает с методом главного критерия, но предлагает некоторый механизм выбора ограничителей, при котором допустимая область в задаче (1.24) гарантированно не пуста.

Существует несколько версий метода уступок. Приведем одну из них, связав ее описание с методом главного критерия. Перенумеруем критерии так, чтобы j -й, используемый в (1.24), стал последним. Выберем величины уступок $\Delta_i \geq 0$, ($i=1, \dots, n$) и выполним следующий итерационный процесс. Вначале положим $Y_0 = Y$ и решим задачу

$$f_1^*(Y_0) = \inf \{ f_1(y) : y \in Y_0 \}.$$

Затем для $i=2, \dots, n$ получим решения серии дополнительных задач

$$f_i^*(Y_{i-1}) = \inf \{ f_i(y) : y \in Y_{i-1}(\Delta_{i-1}) \},$$

где

$$Y_{i-1} = Y_{i-1}(\Delta_{i-1}) = \{ y \in Y_{i-2} : f_{i-1}(y) \leq f_{i-1}^*(Y_{i-2}) + \Delta_{i-1} \}.$$

Если нижняя грань в последней задаче достигается, то точки ее решения, очевидно, являются слабо эффективными для исходной задачи. Это связано с тем, что последняя задача серии эквивалентна методу (1.24) при $j=n$ и ограничителях вида $f_i^+ = f_i^*(Y_{i-1}) + \Delta_i$, обеспечивающих, по построению, непустоту допустимой области в (1.24). То есть за счет выбора $\Delta_i \geq 0$ назначаются уступки по значениям соответствующих критериев f_i ($i \neq j, i=1, \dots, n$).

1.3.5. Метод идеальной точки Вержбицкого

Идея метода заключается в том, чтобы выделить из Y ту точку \underline{y} , f -образ которой наиболее близок к некоторой заданной идеальной оценке \bar{z} в смысле специально введенной функции «расстояния».

Метод формализуется следующим образом. Обозначим через $R^+(z)$ множество векторов z из R^n , с компонентами, большими или равными компонентам z'

$$R^+(z) = \{ z' \in R^n : z' \geq z \}$$

Введем также множество

$$F^+ = \bigcup_{z \in F} R^+(z).$$

Его вид иллюстрируется на рис.1.11.

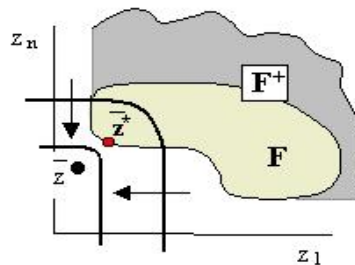


Рис. 1.11. Изолинии функции свертки в методе идеальной точки

Выберем точку $\bar{z} \notin \text{int } F^+$ (через int обозначим внутренность множества), являющуюся параметром метода. Введем скалярную функцию $\Psi_{\bar{z}}(y)$, определяющую «расстояние» в пространстве оценок между f -образом точки y и оценкой \bar{z} . Следует обратить внимание, что данная функция не удовлетворяет всем аксиомам для метрики.

$$\Psi_{\bar{z}}(y) = \rho(f(y), \bar{z}) = \sum_{i=1}^n (\max \{ f_i(y) - \bar{z}_i; 0 \})^2 = \sum_{i=1}^n ((f_i(y) - \bar{z}_i)_+)^2 \quad (1.25)$$

Здесь использовано обозначение

$$(u)_+ = \begin{cases} u, & u \geq 0 \\ 0, & u < 0. \end{cases} \quad (1.26)$$

Далее исходная многокритериальная задача заменяется скалярной задачей вида

$$\Psi_{\bar{z}}(y) = \rho(f(y), \bar{z}) \rightarrow \min, y \in Y. \quad (1.27)$$

Свойство 1.3.

Пусть $y_{\bar{z}}^*$ – решение задачи (1.27) для $\bar{z} \notin \text{int } F^+$, тогда $y_{\bar{z}}^*$ – слабо эффективная точка исходной задачи.


ДОКАЗАТЕЛЬСТВО. Обозначим $\bar{z}^* = f(y_{\bar{z}}^*)$. Пусть свойство неверно, т.е. $\bar{z}^* \notin S(F)$. Отсюда следует, что $\exists y \in Y$, что для $z = f(y)$ выполнится неравенство $z < \bar{z}^*$, т.е. $\forall i = 1, \dots, n \ z_i < \bar{z}_i^*$, откуда

$$z_i - \bar{z}_i < \bar{z}_i^* - \bar{z}_i \quad (1.28)$$

а следовательно, $(z_i - \bar{z}_i)_+ \leq (\bar{z}_i^* - \bar{z}_i)_+$. Равенство возможно для тех i , где разности отрицательны, в силу (1.26). Возводя в квадраты и суммируя, получим, что $\Psi_{\bar{z}}(y) \leq \Psi_{\bar{z}}(y_{\bar{z}}^*)$, а следовательно, в силу оптимальности \bar{z}^* , значения $\Psi_{\bar{z}}(y)$ и $\Psi_{\bar{z}}(y_{\bar{z}}^*)$ совпадают. При этом возможно два случая: оба значения равны нулю или же оба больше нуля.

Рассмотрим первый случай. Если оба значения равны нулю, то тождественно по i будет выполняться $(\bar{z}_i^* - \bar{z}_i)_+ \equiv 0$, т.е. $\bar{z}^* \leq \bar{z}$. Но у нас $z < \bar{z}^*$ и поэтому нашлось $z \in F$, что $z < \bar{z}^* \leq \bar{z}$. Это значит, что $\bar{z} \in \text{int } R^+(z) \subseteq \text{int } F^+$, что противоречит условию.

Рассмотрим второй случай, когда $\Psi_{\bar{z}}(z) > 0$. Тогда хотя бы для одного i $(z_i - \bar{z}_i)_+ > 0$. Для этого значения i будет выполнено строгое неравенство $(z_i - \bar{z}_i)_+ < (\bar{z}_i^* - \bar{z}_i)_+$, в силу специфики операции (1.26) и исходного строгого неравенства (1.28). Отсюда видим, что $\Psi_{\bar{z}}(y) < \Psi_{\bar{z}}(y_{\bar{z}}^*)$, вопреки ранее обнаруженному равенству этих значений. Полученные противоречия доказывают справедливость свойства.

 **Замечание.** Поскольку для $\bar{z} \in S(F)$ среди решений задачи (1.27) обязательно найдутся слабо эффективные точки y^0 , являющиеся прообразами оценки z , то

$$f \left(\bigcup_{\bar{z} \in \text{int } F^+} \{y_{\bar{z}}^*\} \right) = S(F),$$

т.е. за счет подбора идеальной точки \bar{z} можно найти любую слабо эффективную точку среди множества решений задачи (1.25)-(1.27).

Следует отметить, что в методе идеальной точки многокритериальная задача свелась ко множеству однокритериальных задач с тем же допустимым множеством Y , что и исходная. Функция (1.25) обеспечивает свертку векторного критерия в скалярный. Компоненты вектора \bar{z} являются параметрами свертки. Данный метод относится к классу методов свертки.

Рассмотрим еще два вида свертки.

1.3.6. Метод линейной свертки

В этом методе функция свертки линейна

$$\Psi_{\lambda}(y) = \sum_{i=1}^n \lambda_i f_i(y), \quad \lambda_i \geq 0, \quad \sum_{i=1}^n \lambda_i = 1, \quad (1.29)$$

а вспомогательные задачи имеют вид

$$\Psi_{\lambda}(y) \rightarrow \min, \quad y \in Y. \quad (1.30)$$

Решение задачи (1.30) можно интерпретировать геометрически как смещение гиперплоскости $(\lambda, z) = C$ в направлении убывания C при сохранении общих точек этой гиперплоскости со множеством F (рис.1.12).

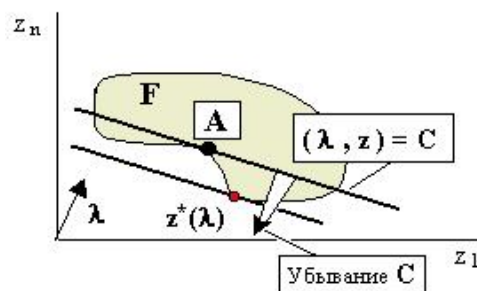


Рис. 1.12. Поиск решения в методе линейной свертки

Свойство 1.4. Пусть y_{λ}^* - решение задачи (1.29)-(1.30), а Y_{λ}^* — множество всех ее глобально-оптимальных решений, тогда справедливо следующее.

- A. В общем случае $y_{\lambda}^* \in Y^0$, а при $\lambda > 0$ $y_{\lambda}^* \in Y^*$.
- B. Если Y выпукло и $f_i(y)$ для $i=1, \dots, n$ – выпуклы, то $\forall z^* \in P(F) \exists \lambda$, что $z^* \in f(Y_{\lambda}^*)$.
- C. В задаче общего вида может найтись оценка $z^* \in P(F)$, что $\forall \lambda \geq 0$ $z^* \notin f(Y_{\lambda}^*)$.

ДОКАЗАТЕЛЬСТВО. Пусть $\lambda > 0$, а $y_{\lambda}^* \notin Y^*$. Тогда $\exists y \in Y$, что $f(y) \neq f(y_{\lambda}^*)$ и $f(y) \leq f(y_{\lambda}^*)$, т.е. $\forall i=1, \dots, n: f_i(y) \leq f_i(y_{\lambda}^*)$ и при этом $\exists j: f_j(y) < f_j(y_{\lambda}^*)$. За счет того, что $\lambda > 0$, получим, что $\Psi_{\lambda}(y) < \Psi_{\lambda}(y_{\lambda}^*)$, но это противоречит оптимальности y_{λ}^* в задаче (1.30).

Если же $\lambda \geq 0$ и $\sum_{i=1}^m \lambda_i = 1$, то в предположении $y_{\lambda}^* \notin Y^0$, получим, что $\exists y \in Y$ с $f_i(y) < f_i(y_{\lambda}^*) \quad \forall i=1, \dots, n$, откуда вновь будет вытекать невозможное неравенство $\Psi_{\lambda}(y) < \Psi_{\lambda}(y_{\lambda}^*)$. Таким образом, свойство (A) верно.

Свойство (B) следует из существования опорной гиперплоскости для выпуклого множества F в любой его граничной точке, а значит и в точках $z^* \in P(F)$. Эта гиперплоскость обеспечивает нестрогую отделимость z^* от F . В качестве вектора λ можно выбрать вектор нормали к этой гиперплоскости. Нестрогий характер разделения приводит к тому, что не обязательно вектор $f(y_{\lambda}^*)$ совпадет с z^* , но для них обязательно совпадут значения функции свертки.

Свойство (C) подтверждается примером. Достаточно выбрать z^* , совпадающей с точкой A на рис.1.12.

Последнее свойство показывает, что метод линейной свертки в общем случае не позволяет выделить все эффективные решения.

1.3.7. Свертка Ю.Б. Гермейера

Очевидно, что для возможности выделения всех эффективных решений необходимо изменить вид поверхности равного уровня функций свертки в пространстве критериев. Свертка, предложенная Ю.Б. Гермейером имеет поверхности равного уровня в виде «угла» (см. рис.1.13), где под «углом» понимается граница множества вида, $R^-(z') = \{z \in R^n : z < z'\}$. Эта свертка определяется соотношением

$$\Psi_\lambda(y) = \max \{ \lambda_i f_i(y) : i = 1, \dots, n \}, \lambda_i \geq 0 \sum_{i=1}^n \lambda_i = 1 \quad (1.31)$$

Замечание. Свертку Гермейера используют в том случае, когда $F \subset R^+(0)$, т.е. когда $\forall z \in F z > 0$. Выполнения этого требования всегда можно добиться, сделав замену переменных в пространстве оценок $z := z + f^*$. При этом придется предварительно решить n вспомогательных задач вида

$$f_i^* = \min \{ f_i(y) : y \in Y \}.$$

По отношению к свертке (1.31) решается задача (1.30).

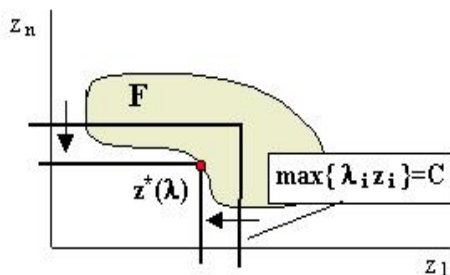


Рис. 1. 13. Изолинии свертки Ю.Б. Гермейера

Свойство 1.5. Если $F \subset R^+(0)$, то $\forall y^o \in Y^o$ найдется такое $\lambda \geq 0$, $\sum_{i=1}^n \lambda_i = 1$, что $\min \{ \Psi_\lambda(y) : y \in Y \} = \Psi_\lambda(y^o)$, т.е. любое слабо эффективное решение может быть получено из решения задачи (1.31)–(1.30) при соответствующем выборе λ .

ДОКАЗАТЕЛЬСТВО проиллюстрируем геометрически. Пусть $z^o = f(y^o)$. Подберем λ так, чтобы в пространстве переменных z поверхность уровня свертки Гермейера

$$\max \{ \lambda_i z_i : i = 1, \dots, n \} = \max \{ \lambda_i f_i(y^o) : i = 1, \dots, n \}$$

имела угловую точку (см. рис.1.13) именно в точке $f(y^o)$. Это определит нужные значения λ .

Теперь докажем формально. Выберем $\lambda_i = (1/z_i^o) / \sum_{j=1}^n (1/z_j^o)$. Предположим, что этот выбор не верен, тогда при таких λ_i найдется $z \in F$, что

$$\max \{ \lambda_i z_i : i = 1, \dots, n \} < \max \{ \lambda_i z_i^o : i = 1, \dots, n \} = \left(\sum_{j=1}^n (1/z_j^o) \right)^{-1}.$$

Подставляя выражение для λ_i и умножая полученное неравенство на постоянную в его конце сумму, получим, $\max \left\{ \frac{z_i}{z_i^0} : i = 1, \dots, n \right\} < 1$, т.е. $z_i < z_i^0 \quad \forall i$, что противоречит исходному предположению об оптимальности по Слейтеру оценки $z^0 = f(y^0)$. Свойство доказано.

1.3.8 Проблема оценивания всего множества эффективных точек.

Приведенные в пунктах 1.3.2–1.3.6 схемы компромисса позволяют выделять отдельные эффективные или слабо эффективные точки, сводя задачу их оценивания к однокритериальным (скалярным) задачам математического программирования. Изменяя значения параметров выбранной схемы компромисса (величин уступок, идеальную точку z , параметры сверток λ) можно получать разные решения многокритериальной задачи. Однако остается открытым вопрос о построении

ε -покрытий множеств Y^* или Y^0 , поскольку в рамках изложенных подходов не ясно как надо изменять параметры схем компромисса для получения подобного покрытия. Следовательно, нужны иные подходы.

Конкретизируем понятие ε -покрытия на примере множества эффективных точек Y^* (назовем такое покрытие ε -оптимальным решением).

Определение 1.6. Множество $Y_\varepsilon^* \subseteq Y^*$ ($\varepsilon > 0$) будем называть ε -оптимальным (по Парето) решением задачи, если $\forall y^* \in Y^* \exists y_\varepsilon^* \in Y_\varepsilon^*$, что $f(y_\varepsilon^*) \leq f(y^*) + \varepsilon \cdot e$, где $e = (1, \dots, 1)^T$ и в Y_ε^* нет двух разных точек $y_\varepsilon^{*1}, y_\varepsilon^{*2}$, что $f(y_\varepsilon^{*1}) \leq f(y_\varepsilon^{*2})$.

Позднее в главе 6 будут рассмотрены специальные методы построения ε -оптимальных решений многокритериальных задач.

1.4. Модели функций, используемые в задачах оптимального выбора

Методы отыскания решений в задачах оптимального выбора можно разделить на две большие группы. К первой относятся методы аналитического и численно-аналитического решения, основанные на использовании условий экстремума (Лагранжа, Куна-Таккера и их обобщений). Ко второй – методы вычислительной оптимизации, использующие измерения локальных характеристик целевых функций и функций ограничений (т.е. значений этих функций, а также, возможно, их градиентов и матриц Гессе) для организации процесса поиска решения (такие измерения в дальнейшем будем называть *испытаниями*).

Совокупность результатов всех проведенных к k -му шагу испытаний функции f называют *поисковой информацией* по этой функции. В дальнейшем будем использовать для нее обозначение $\omega_k = \omega(f, Y_k)$, где $Y_k = \{y^1, \dots, y^k\}$.

При любом способе поиска решения необходима предварительная (*априорная*) информация о свойствах решаемой задачи. В первом случае, при аналитическом решении задачи, на основе априорной информации выбирается конкретный вид условий экстремума, а во втором, при использовании вычислительных методов, априорная информация позволяет интерпретировать

результаты проведенных испытаний с целью планирования следующих измерений и, тем самым, определяет выбор вычислительного метода.

Хорошей иллюстрацией последнего тезиса является задача поиска минимума функции $f(y)$ на отрезке $y \in D=[a,b]$. Пусть проведено $k=2$ измерений функции в точках $a < y^1 < y^2 < b$, причем оказалось, что $f^1 < f^2$. Если априори известно, что функция f принадлежит классу функций $\Phi = U[D]$, унимодальных на отрезке D , то решение y^* может находиться только в подобласти $D_k = [a, y^2]$, и дальнейшие испытания следует проводить только в ней, если же дополнительно известно, что функция f является липшицевой на $D = [a,b]$ с константой L , т.е. $f \in \Phi = U[D] \cap Lip[D]$, то при тех же результатах измерений решение y^* принадлежит уже области D_k другого вида, а именно $D_k = [a, y^2 - (f^2 - f^1)/L]$, следовательно, получение дополнительной информации приводит к изменению областей размещения последующих испытаний (рис.1.14).

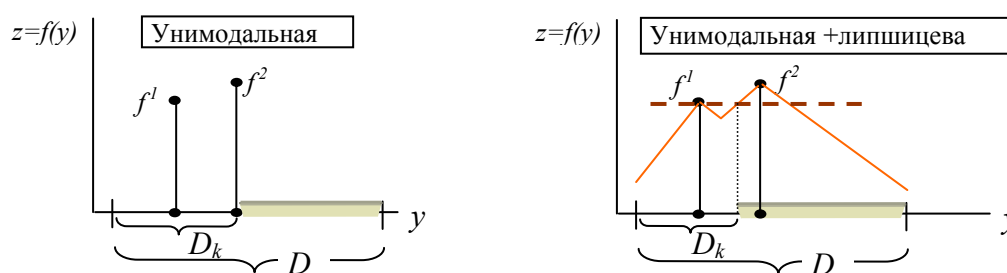


Рис.1.14 Различное сокращение области поиска в зависимости от априорной информации

Определение 1.7. Под моделью задачи оптимального выбора будем понимать совокупность постановки задачи вместе с имеющимися предположениями о свойствах функций, входящих в постановку.

Методы решения всегда выбираются из принятой модели задачи. В дальнейшем будем употреблять также термины «модель целевой функции», «модель функций ограничений», понимая под этим совокупность предположений о свойствах этих функций.

В этом разделе будут рассмотрены модели функций, наиболее часто используемые при построении и исследовании вычислительных методов оптимизации, а так же в теории условий экстремума. Первая группа моделей особенно важна для теории локальной оптимизации и условий экстремума. Это выпуклые функции и их обобщения. Вторая группа моделей важна для многоэкстремальной оптимизации. Это Липшицевы, Гельдеровы функции, функции с липшицевыми производными по направлениям.

1.4.1. Модели функций, основанные на представлениях о выпуклости

Известно несколько подходов к понятию выпуклости скалярной функции $f: Y \rightarrow R^1$, определенной на выпуклом множестве Y .

1.4.1.1. Выпуклые, строго и сильно выпуклые функции

Материал этого пункта необходим при изучении второй и седьмой глав.

Обычное понятие *выпуклости функции* вводится как свойство, эквивалентное требованию выпуклости множества, называемого *надграфиком* или *эпиграфом*.

Определение 1.8. Надграфиком (эпиграфом) функции $f(y)$ называется множество

$$\text{epi } f = \{(y, z) : y \in Y \subseteq \mathbb{R}^N, z \in \mathbb{R}^1, z \geq f(y)\} \quad (1.32)$$

Определение 1.9. Функция f , определенная на выпуклом множестве Y , называется выпуклой (вогнутой), если

$$\forall y^1, y^2 \in Y \text{ и } \forall \alpha \in (0,1) f(\alpha y^1 + (1-\alpha)y^2) \leq (\geq) \alpha f(y^1) + (1-\alpha)f(y^2) \quad (1.33)$$

Понятие строгой выпуклости (строгой вогнутости) соответствуют случаю строгого выполнения неравенства.

Оказывается, что выпуклость функции равносильна выпуклости ее надграфика (рис. 1.15).

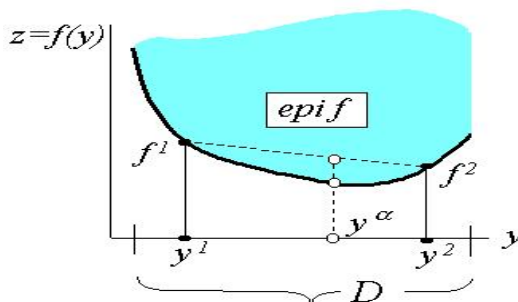


Рис.1.15. Выпуклость функции и ее надграфика

Теорема 1.1. Пусть Y -выпукло. Функция $f(y)$, выпукла на Y тогда и только тогда, когда выпукло множество $\text{epi } f$.

НЕОБХОДИМОСТЬ. Для произвольных точек (y^1, z^1) и (y^2, z^2) из $\text{epi } f$ выполняется неравенство $z^i \geq f(y^i)$ ($i=1,2$).

Пусть $(y^\alpha, z^\alpha) = \alpha(y^1, z^1) + (1-\alpha)(y^2, z^2)$. Тогда из выпуклости $f(y)$ для $\alpha \in (0,1)$ вытекает неравенство $f(y^\alpha) \leq \alpha f(y^1) + (1-\alpha)f(y^2) \leq \alpha z^1 + (1-\alpha)z^2 = z^\alpha$, которое означает принадлежность точки (y^α, z^α) надграфу $\text{epi } f$.

ДОСТАТОЧНОСТЬ. Легко видеть, что для любых y^1, y^2 из Y точки $(y^i, f(y^i)) \in \text{epi } f$ для $i = 1,2$. Из выпуклости надграфика при любом $\alpha \in (0,1)$ следует, что точка

$$(y^\alpha, z^\alpha) = (\alpha y^1 + (1-\alpha)y^2, \alpha f(y^1) + (1-\alpha)f(y^2)) \in \text{epi } f,$$

что означает выполнение определения выпуклости.

Выпуклые функции обладают рядом полезных свойств. Укажем на некоторые из них.

Свойство 1.6. Множества $Y_C = \{y \in Y : f(y) \leq C\}$ выпуклы.

Свойство 1.7. Пусть Y – не пусто, выпукло и открыто, f – дифференцируема на Y . При этом f выпукла (строго выпукла) на Y тогда и только тогда, когда $\forall y, y \in Y$

$$f(y) \geq (>) f(\bar{y}) + (\nabla f(\bar{y}), y - \bar{y}). \quad (1.34)$$

Свойство 1.8. Пусть Y – не пусто, выпукло и открыто, $f(y) \in C^2(Y)$. При этом f выпукла на Y тогда и только тогда, когда $\forall y \in Y$ матрица Гессе $\Gamma^f(y)$ неотрицательно определена.

Свойство 1.9. Любой локальный минимум выпуклой функции f на выпуклом множестве Y является глобальным, а само множество глобальных минимумов Y^* выпуклым. У строго выпуклой функции существует единственный локальный минимум $y^0 = y^*$.

Замечание. Свойство 1.8 не имеет точного аналога для строго выпуклых функций. А именно, если $\Gamma^f(y)$ – положительно определена для $y \in Y$, то f – строго выпукла на Y , однако обратное не верно. Из строгой выпуклости f в общем случае следует лишь неотрицательная определенность матриц Гессе.

Справедливость этого замечания подтверждает следующий простой контр пример. Для строго выпуклой в R^1 функции $f(y) = y^4$ значение $\Gamma^f(y) = 12y^2$ и обращается в 0 при $y = 0$.

Следует обратить внимание на содержательную сторону свойств. Первое из них обеспечивает выпуклость множеств точек, ограниченных поверхностями равного уровня функции f . Второе свойство позволяет по измерениям функции и ее градиента в точке \bar{y} отсекал области, не содержащие решения. А именно, можно отсечь множество точек со значениями $(\nabla f(\bar{y}), y - \bar{y}) > 0$, что позволяет сокращать область поиска решения (рис.1.16). Наконец, третье свойство связывает выпуклость дважды непрерывно дифференцируемой функции с неотрицательностью собственных чисел матриц вторых производных функции f .

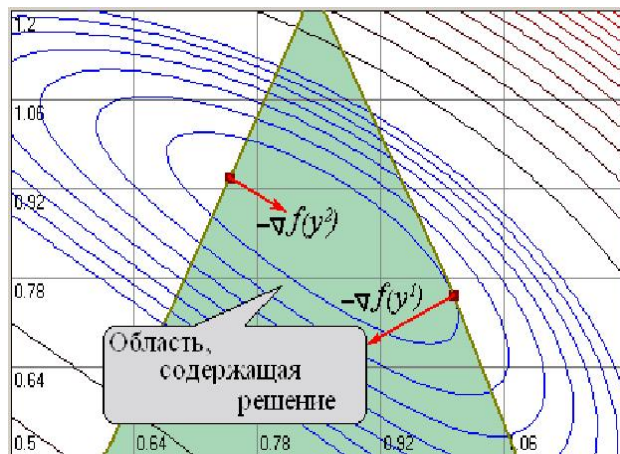


Рис.1.16. Сокращение области поиска для выпуклой дифференцируемой функции

Введем теперь понятие *сильно выпуклой функции*, важное при исследовании вычислительных методов.

Определение 1.10. Функция $f(y)$, определенная на выпуклом множестве Y называется *сильно выпуклой* (с параметром $\rho > 0$), если $\forall y^1, y^2 \in Y$ и $\forall \alpha \in (0, 1)$ выполняется

$$f(\alpha y^1 + (1 - \alpha)y^2) \leq \alpha f(y^1) + (1 - \alpha)f(y^2) - \rho\alpha(1 - \alpha)\|y^1 - y^2\|^2. \quad (1.35)$$

Для функций этого класса выполняются следующие усиленные варианты свойств 1.7, 1.8.

Свойство 1.7'. Пусть f – дифференцируемая функция, определенная на непустом открытом выпуклом множестве Y . f сильно выпукла (с параметром ρ) тогда и только тогда, когда $\forall y, \bar{y} \in Y$ выполняется неравенство

$$f(y) \geq f(\bar{y}) + (\nabla f(\bar{y}), y - \bar{y}) + \rho \|y - \bar{y}\|^2. \quad (1.36)$$

Приведем ДОКАЗАТЕЛЬСТВО НЕОБХОДИМОСТИ. Пусть $y^\alpha = \alpha y + (1 - \alpha)\bar{y}$. Из определения, при $\alpha \in (0, 1)$ следует, что

$$f(y^\alpha) \leq f(\bar{y}) + (f(y) - f(\bar{y}))\alpha - \rho \|y - \bar{y}\|^2 \alpha(1 - \alpha).$$

Отсюда

$$(f(\bar{y} + (y - \bar{y})\alpha) - f(\bar{y})) / \alpha \leq f(y) - f(\bar{y}) - \rho \|y - \bar{y}\|^2 (1 - \alpha).$$

При $\alpha \rightarrow 0$ в пределе получим неравенство (1.4.5).


Свойство 1.8'. Пусть выполнены предположения свойства 1.8. При этом f сильно выпукла (с некоторым параметром $\rho > 0$) тогда и только тогда, когда найдется такое $m > 0$, что $\forall y \in Y$ и $\forall d \neq 0$

$$d^T \Gamma^f(y) d \geq m \|d\|^2, m > 0 \quad (1.37)$$

Приведем ДОКАЗАТЕЛЬСТВО ДОСТАТОЧНОСТИ. Пусть (1.37) выполнено, тогда $\exists \theta \in [y, \bar{y}]$, что

$$\begin{aligned} f(y) &= f(\bar{y}) + (\nabla f(\bar{y}), y - \bar{y}) + 0,5(y - \bar{y})^T \Gamma^f(\theta)(y - \bar{y}) \geq \\ &\geq f(\bar{y}) + (\nabla f(\bar{y}), y - \bar{y}) + 0,5 m \|y - \bar{y}\|^2. \end{aligned}$$

Отсюда по свойству 1.7' функция f будет сильно выпукла с параметром $\rho = 0,5m$.

 **Замечание.** Неравенство (1.37) эквивалентно условию

$$\exists m: 0 < m \leq \lambda_i \quad (i=1, 2, \dots, N) \quad (1.38)$$

где λ_i – собственные числа матрицы $\Gamma^f(y)$.

Это вытекает из того, что матрица Гессе эквивалентна диагональной матрице Λ с числами $\lambda_1, \dots, \lambda_N$ на диагонали: $\Gamma^f(y) = R^T \Lambda R$, $R^T = R^{-1}$. При этом замена $z = Rd$ позволяет представить (1.37) в виде:

$$\forall z \neq 0 \quad \sum_{i=1}^N z_i^2 \lambda_i \geq \sum_{i=1}^N z_i^2 m,$$

откуда следует справедливость замечания.

Таким образом, сильная выпуклость для функции $f \in C^2(Y)$ эквивалентна равномерной (по y) отделенности от нуля собственных чисел матриц Гессе $\Gamma^f(y)$.

Сильно выпуклые функции обладают рядом дополнительных замечательных свойств, позволяющих оценивать ошибку по координате в определении решения. Эти свойства являются следствиями из свойства 1.7'.

Следствие 1.7'.1. В условиях свойства 1.7' для сильной выпуклости функции необходимо и достаточно, чтобы $\forall y, \bar{y} \in Y$ и $\rho > 0$ выполнялось

$$(\nabla f(y) - \nabla f(\bar{y}), y - \bar{y}) \geq 2\rho \|y - \bar{y}\|^2 \quad (1.39)$$

ДОКАЗАТЕЛЬСТВО НЕОБХОДИМОСТИ легко получается, если сложить (1.36) с аналогичным неравенством, в котором y и \bar{y} поменяли местами.

Следствие 1.7'2. В условиях свойства 1.7' для точки y^* – внутреннего глобального минимума сильно выпуклой (с параметром ρ) функции справедливы следующие оценки

$$\rho \|y - y^*\|^2 \leq f(y) - f(y^*) \quad (1.40)$$

$$2\rho \|y - y^*\| \leq \|\nabla f(y)\| \quad (1.41)$$

Эти оценки непосредственно вытекают из (1.36) и (1.39), если принять $\bar{y} = y^*$ и учесть, что $\nabla f(y^*) = 0$.

1.4.1.2. Квазивыпуклые, строго и сильно квазивыпуклые функции

Материал пунктов 1.4.1.2 и 1.4.1.3 является факультативным. При чтении его можно опустить и это существенно не повлияет на понимание следующего материала.

Квазивыпуклость является одним из обобщений понятия выпуклой функции. Дело в том, что для выполнения некоторых полезных свойств выпуклость является излишне жестким требованием. Можно наложить на функцию более мягкие требования, сохранив некоторые нужные свойства.

Введем другое понятие выпуклости (*квазивыпуклость*) эквивалентное выпуклости множеств следующего вида

$$Y_C = \{y \in Y : f(y) \leq C\} \quad (1.42)$$

для любых C (свойство 1.6 выпуклых функций).

Определение 1.11. Функцию f , определенную на пустом выпуклом множестве Y , назовем квазивыпуклой, если $\forall y^1, y^2 \in Y$ и $\forall \alpha \in (0, 1)$

$$f(\alpha y^1 + (1 - \alpha)y^2) \leq \max\{f(y^1); f(y^2)\} \quad (1.43)$$


Теорема 1.2. Функция f квазивыпукла на выпуклом Y тогда и только тогда, когда $\forall C$ множества Y_C из (1.42) выпуклы.

Необходимость вытекает из того, что из условия $y^1, y^2 \in Y_C$ и (1.43) следует, что

$$f(y^\alpha) = f(\alpha y^1 + (1 - \alpha)y^2) \leq \max\{f(y^1); f(y^2)\} \leq C,$$

следовательно, $y^\alpha \in Y_C$.

Достаточность. Пусть (1.42) верно $\forall C$. Для произвольных $y^1, y^2 \in Y_C$ достаточно выбрать $C = \max\{f(y^1), f(y^2)\}$. Тогда $y^1, y^2 \in Y_C$, а в силу его выпуклости $\forall \alpha \in (0, 1)$ $y^\alpha = \alpha y^1 + (1 - \alpha)y^2 \in Y_C$. Это значит, что $f(y^\alpha) \leq C = \max\{f(y^1); f(y^2)\}$, что эквивалентно (1.43).

 **Замечание.** Хорошо известное свойство класса выпуклых функций — замкнутость по отношению к операции сложения. Это свойство для квазивыпуклых функций нарушается, т.е. сумма двух квазивыпуклых функций может не являться квазивыпуклой (например, $f(y_1, y_2) = y_1^3 + y_2^3$).

Примеры поведения квазивыпуклых функций представлены на рис.1.17.

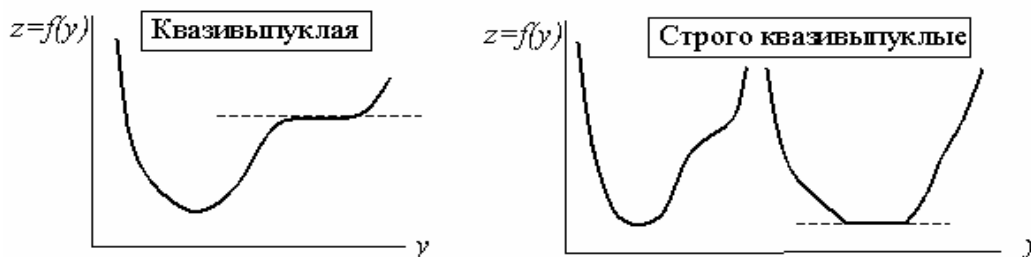


Рис.1.17. Примеры поведения квазивыпуклых функций

Нетрудно видеть, что квазивыпуклость допускает существование у функции зон стационарности со значениями, отличными от глобально оптимального. Это приводит к нарушению важного свойства 1.9, характерного для выпуклых функций.

Для того чтобы это свойство было выполнено, используют понятие *строгой квазивыпуклости* (см. рис.1.17).

Определение 1.12. Функция f , определенная на непустом выпуклом множестве, называется *строго квазивыпуклой*, если $\forall y^1, y^2 \in Y$, что $f(y^1) \neq f(y^2) \quad \forall \alpha \in (0, 1)$ выполняется

$$f(\alpha y^1 + (1 - \alpha)y^2) < \max \{f(y^1); f(y^2)\} \quad (1.44)$$

Свойство 1.9'а. Для строго квазивыпуклой функции любой локальный минимум является глобальным.

ДОКАЗАТЕЛЬСТВО. Если предположить существование локального y^o и глобального y^* минимумов со значениями $f(y^*) < f(y^o)$, то мы приходим к противоречию, поскольку при $\alpha \rightarrow 0$ точка $y^\alpha = \alpha y^* + (1 - \alpha)y^o \rightarrow y^o$, но $f(y^\alpha) < \max \{f(y^*); f(y^o)\} = f(y^o)$, и y^o не может быть локальным минимумом.

К сожалению, определение 1.12 имеет существенный недостаток, связанный с тем, что строго квазивыпуклая функция может не являться квазивыпуклой. Это связано с тем, что в определении ничего не сказано о свойствах функции при $f(y^1) = f(y^2)$. Существует еще одно понятие — *квазивыпуклость*, исправляющее эту ситуацию.

Определение 1.13. Функция $f(y)$, определенная на непустом выпуклом Y , называется *сильно квазивыпуклой*, если $\forall y^1, y^2 \in Y$ и $y^1 \neq y^2 \quad \forall \alpha \in (0, 1)$ выполняется (1.44).

Свойство 1.9'в. У сильно квазивыпуклой функции существует единственный локальный минимум, совпадающий с глобальным.

ДОКАЗАТЕЛЬСТВО. В силу предыдущего свойства все локальные минимумы будут глобальными. Далее, если бы $\exists y^1 \neq y^2$ и $f(y^1) = f(y^2) = f(y^*)$ то, $\forall \alpha \quad f(y^\alpha) < \min \{f(y^1); f(y^2)\} = f(y^*)$, что является противоречием к локальной оптимальности точек y^1, y^2 .

Вновь вернемся к квазивыпуклым функциям. Для них выполняется аналог свойства 1.7 выпуклых функций.

Свойство 1.7". Пусть $f(y)$ определена и дифференцируема на непустом, выпуклом множестве Y . При этом f квазивыпукла на Y тогда и только тогда, когда $\forall y, \bar{y} \in Y$ выполнено любое из двух эквивалентных условий

A. Если $f(y) \leq f(\bar{y})$, то $(\nabla f(\bar{y}), y - \bar{y}) \leq 0$ (1.45)

B. Если $(\nabla f(\bar{y}), y - \bar{y}) > 0$, то $f(y) > f(\bar{y})$, (1.46)

ДОКАЗАТЕЛЬСТВО. Эквивалентность условий A и B непосредственно доказывается методом от противного. Докажем необходимость условия A. Пусть f – квазивыпукла и $f(y) \leq f(\bar{y})$. Воспользуемся разложением

$$f(\alpha y + (1-\alpha)\bar{y}) - f(\bar{y}) = f(\bar{y} + \alpha(y - \bar{y})) - f(\bar{y}) = (\nabla f(\bar{y}), y - \bar{y})\alpha + o(\alpha).$$

Из квазивыпуклости f при $\alpha \in (0, 1)$

$$f(\alpha y + (1-\alpha)\bar{y}) \leq \max \{f(y); f(\bar{y})\} = f(\bar{y}),$$

что говорит о неположительности последней суммы в предыдущем равенстве. Устремляя α к нулю, получаем $(\nabla f(\bar{y}), y - \bar{y}) \leq 0$.


Доказательство достаточности можно найти в книге [1].

Свойства A и B очевидно выполняются и для сильно квазивыпуклых функций, поэтому свойство 1.7" можно использовать для них, чтобы проводить отсечения областей не содержащих решения.

A именно, множества вида

$$Y(\bar{y}) = \{y \in Y : (\nabla f(\bar{y}), y - \bar{y}) > 0\} \quad (1.47)$$

не содержит точек глобального минимума функции f на Y .

 **Замечание.** Заметим, что квазивыпуклые и сильно квазивыпуклые модели поведения функций обычно используются в тех ситуациях, когда требуется обеспечить выпуклость областей, выделяемых ограничениями равенствами и неравенствами при возможно более слабых предположениях о функциях ограничений. Вторым случаем использования этих моделей является возможность применять правила $y^* \notin Y(\bar{y})$, вытекающие из (1.47) при более слабых предположениях о функции f , чем обычная выпуклость.

1.4.1.3. Псевдовыпуклые и строго псевдовыпуклые функции

Понятие псевдовыпуклости вводится только для дифференцируемых функций. Его определяют таким образом, чтобы из равенства $\nabla f(y^0) = 0$, $y^0 \in Y$ следовала глобальная оптимальность точки y^0 на открытом множестве Y . В качестве определения используются условия, близкие к A и B из свойства 1.7".

Определение 1.4.8. Пусть f – дифференцируема на непустом открытом выпуклом множестве Y . Говорят, что $f(y)$ псевдовыпукла (строго псевдовыпукла) если $\forall y, \bar{y} \in Y$ выполняется одно из двух эквивалентных условий A' или B':

A'. Если $\nabla f(y) < (\leq) f(\bar{y})$, то $(\nabla f(\bar{y}), y - \bar{y}) < 0$ (1.48)

B'. Если $(\nabla f(\bar{y}), y - \bar{y}) \geq 0$, то $f(y) \geq (>) f(\bar{y})$ (1.49)

На рис.1.18 приведен пример поведения псевдовыпуклой функции.

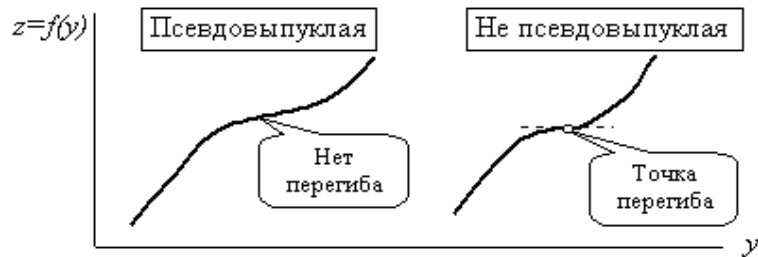


Рис.1.18. Примеры поведения псевдовыпуклых функций

Приведенное ниже свойство можно считать аналогом свойства 1.9 для выпуклых функций.

Свойство 1.9'. Если f – псевдовыпукла на открытом Y , то из равенства $\nabla f(y^0)=0$, $y^0 \in Y$ следует, что y^0 – глобальный минимум f на Y .

ДОКАЗАТЕЛЬСТВО. Поскольку $\nabla f(y^0)=0$, то $\forall y \in O_\varepsilon(y^0)$ $(\nabla f(y^0), y-y^0)=0$. Тогда по условию $B' f(y) \geq f(y^0)$, т.е. y^0 -локальный минимум. Ниже мы покажем, что из псевдовыпуклости вытекает строгая квазिवыпуклость, поэтому локальный минимум является глобальным по свойству 1.9'а.

Необходимая нам взаимосвязь между несколькими типами выпуклости устанавливается следующей теоремой.

Теорема 1.3. Если f – псевдовыпукла на Y , то f строго квазивыпукла и квазивыпукла.

ДОКАЗАТЕЛЬСТВО. Возьмем $y^1, y^2 \in Y$, где $f(y^1) < f(y^2)$. Допустим, строгой квазивыпуклости нет, тогда возможна ситуация, когда для $\alpha \in (0,1)$ и $y^\alpha = \alpha y^1 + (1-\alpha)y^2$

$$f(y^\alpha) \geq \max\{f(y^1), f(y^2)\} = f(y^2). \quad (1.50)$$

Из псевдовыпуклости, неравенство $f(y^1) < f(y^2) \leq f(y^\alpha)$ влечет $(\nabla f(y^\alpha), y^1 - y^\alpha) < 0$. Поскольку $y^1 - y^\alpha = -\beta(y^2 - y^\alpha)$ при некотором $\beta > 0$, то

$$(\nabla f(y^\alpha), y^2 - y^\alpha) > 0 \quad (1.51)$$

и, следовательно, из условия B' получаем, $f(y^2) \geq f(y^\alpha)$. Но тогда, учитывая обратное неравенство (1.50), получаем $f(y^2) = f(y^\alpha)$.

Вернемся к неравенству (1.51). Оно означает, что на интервале $(y^\alpha; y^2)$ найдется точка \bar{y} , где $f(\bar{y}) > f(y^\alpha) = f(y^2)$. Применим дважды условие A' , получим неравенства

$$(\nabla f(\bar{y}), y^2 - \bar{y}) < 0 \text{ и } (\nabla f(\bar{y}), y^\alpha - \bar{y}) < 0,$$

но это невозможно, поскольку $\bar{y} \in (y^\alpha, y^2)$ и существует $\gamma > 0$, что $y^2 - \bar{y} = -\gamma(y^\alpha - \bar{y})$. Итак, f строго квазивыпукла.

Для доказательства квазивыпуклости необходимо дополнительно показать, что при $f(y^1) = f(y^2)$ и $y^1 \neq y^2$ $f(y^\alpha) \leq \max\{f(y^1); f(y^2)\} = f(y^1) = f(y^2)$. Если это не так, то существует $\alpha \in (0,1)$, что $f(y^\alpha) > f(y^1) = f(y^2)$. Тогда из условия A' получаем $(\nabla f(y^\alpha), y^1 - y^\alpha) < 0$ и $(\nabla f(y^\alpha), y^2 - y^\alpha) < 0$, что противоречит неравенству (1.51).

1.4.2. Модели функций используемые в многоэкстремальной оптимизации

Материал этого пункта будет широко использоваться в главах 4, 5 и 6.

Задачи многоэкстремальной оптимизации принципиально отличаются от более простых одноэкстремальных задач. Причина кроется в кардинальном отличии понятий глобального (абсолютного) и локального экстремумов. Глобальный экстремум является интегральной характеристикой задачи и его отыскание, с заданной точностью, требует построения оценок поведения функций задачи в целом во всей области поиска. При этом следует иметь ввиду, что в зависимости от характера имеющейся априорной информации о свойствах решаемой задачи эти оценки могут быть двух типов: гарантирующие (детерминированные) и вероятностные. В соответствии с этим мы будем говорить как о гарантирующих (детерминированных) моделях поведения функций, так и вероятностных моделях их поведения.

В рамках детерминированных моделей по конечному числу испытаний функций задачи должно быть возможно построение оценок их значений в точках области определения, а также оценок положения решений рассматриваемой задачи оптимального выбора.

Вероятностные модели представляют функции задачи как реализации случайного процесса (или поля) с заданными свойствами. Учет поисковой информации обычно осуществляется переходом к их апостериорному вероятностному описанию, условному по отношению к дополнительной полученной информации. При этом оценки поведения функций и оценки решения носят вероятностный характер. В последующих главах будет показано, каким образом в рамках принятых моделей поведения функций можно строить алгоритмы поиска оптимальных решений, как оптимальные решающие правила.

1.4.2.1. Примеры детерминированных моделей многоэкстремальных функций

Наиболее простой и часто используемой моделью данного типа является липшицева модель. Она позволяет учитывать поисковую информацию по вычисленным значениям функций и адекватна тем задачам, в которых функции имеют ограниченные скорости изменения при варьировании переменных.

Определение 1.15. Будем говорить, что функция $f(y)$ липшицева в области D с константой L , т.е. $f \in \Phi = Lip(D)$, если $\forall y^1, y^2 \in D$

$$|f(y^1) - f(y^2)| \leq L \|y^1 - y^2\| \quad (1.52)$$

Очевидно, для липшицевых функций по конечному числу измерений y^i , $f^i = f(y^i)$ ($i=1, \dots, k$) можно построить следующие оценки

$$\forall y \in D: f_k^-(y) \leq f(y) \leq f_k^+(y),$$

$$f_k^+(y) = \min \{f^i + L \|y - y^i\| : i = 1, \dots, k\} \quad (1.53)$$

$$f_k^-(y) = \max \{f^i - L \|y - y^i\| : i = 1, \dots, k\} \quad (1.54)$$

Если ввести специальное обозначение для достигнутого минимального значения функции $f(y)$

$$f_k^* = \min \{f^i : i = 1, \dots, k\},$$

то используя (1.54), легко записать оценки глобального минимума y^* функции f на D как по координатам, так и по значению функции:

$$f_k^- = \min\{f_k^-(y) : y \in D\} \leq f(y^*) \leq f_k^* \quad (1.55)$$

$$y^* \in D_k = \{y \in D : f_k^-(y) \leq f_k^*\} \quad (1.56)$$

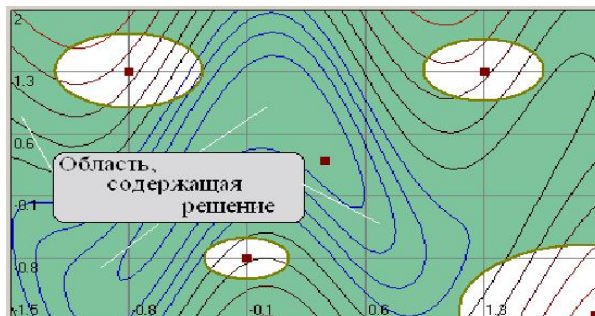


Рис.1.19. Сокращение области поиска при измерениях липшицевой функции двух переменных

Соотношение (1.55) позволяет определять погрешность оценки f_k^* по значению функции, а свойство (1.56) характеризует области целесообразного размещения следующих испытаний (см. пример для задачи с $y \in D \subset R^2$ при $k=5$ на рис.1.19).

В случае одного переменного приведенные свойства иллюстрируются на рис.1.20.

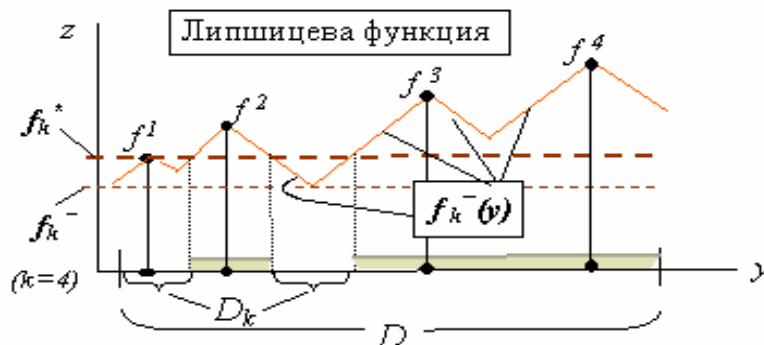


Рис.1.20. Сокращение области поиска при измерениях липшицевой функции от одной переменной

Пример соответствует тому случаю, когда в R^1 $\|y\| = |y|$. Заметим, что если в R^1 выбирать $\|y\| = |y|^{1/p}$, то функцию f , липшицеву в этой метрике, обычно называют *гельдеровой*, а условие (1.52) — *условием Гельдера*. На рис.1.21 показан соответствующий вид оценок по Гельдеру при $p = 2$ и одинаковых с рис.1.20 значениях константы L . Следует заметить, что при достаточной близости точек измерений y^i условие Гельдера приводит к более медленному сокращению меры области поиска D_k по сравнению с условием Липшица.

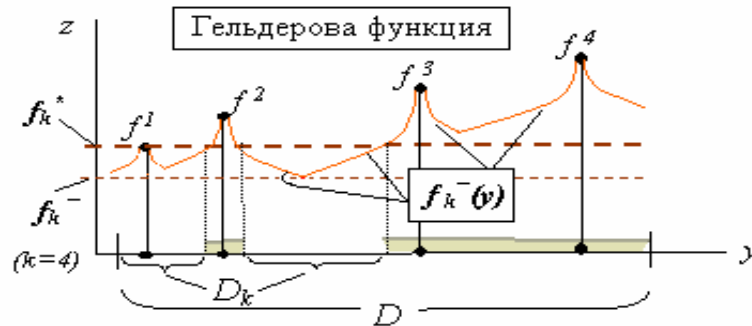


Рис.1.21. Сокращение области поиска при измерениях гелдеровой функции

Первый вычислительный метод оптимизации, использующий липшицеву модель был построен Пиявским С.А. [37] в 1967 году. Подробнее этот материал будет представлен в главах 4, 5.

Замечание. Существенный недостаток модели типа (1.52) заключается в том, что константа L принимается одинаковой для всей области D , тогда как для повышения точности оценок функции (1.53), (1.54) следует считать значение L зависящим от положения точки y в D (L обычно убывает при приближении точки y к решению). При построении вычислительных методов на основе моделей типа (1.52) такая зависимость может быть учтена.

1.4.2.2. Примеры вероятностных моделей многоэкстремальных функций

Вероятностная модель функции f предполагает задание вероятностной меры P на системе (измерениях) подмножеств некоторого класса функций Φ , заданных на D . Тем самым неизвестная многоэкстремальная функция может быть рассмотрена как реализация случайного поля, а в одномерном случае ($D \subseteq R^1$) — случайного процесса. Такой подход позволяет применить к построению процедур поиска методы теории оптимальных статистических решений. Подробнее этот вопрос будет рассмотрен в главе 4.

Сейчас следует отметить, что в 60-х годах были предложены два существенно различных способа описания и использования вероятностных моделей в многоэкстремальной оптимизации. При этом авторы работ не знали о работах друг друга. Один вариант был предложен Кушнером Х. [50, 51], а другой — в совместной работе Неймарком Ю.И. и Стронгиным Р.Г. [35].

В этом пункте мы поясним возможный принцип построения вероятностной модели на примере функции одной переменной $f(y)$, $y \in D = [a, b] \subset R^1$, следуя идеям Кушнера Х.

В его модели функция $f(y)$ считается непрерывной реализацией гауссова случайного процесса $\xi(y)$ с независимыми приращениями (винеровского процесса). Это значит, что для любых $y^1 \leq y^2 \leq y^3 \leq y^4, \dots$ случайные величины приращений $\xi(y^2) - \xi(y^1)$, $\xi(y^3) - \xi(y^2), \dots$ являются независимыми и любое из этих приращений распределено по нормальному закону со средним 0 и дисперсией, пропорциональной приращению аргумента, т.е. разности $\xi(y^{i+1}) - \xi(y^i)$ распределены по нормальному закону $N(0, \sigma / y^{i+1} - y^i /)$.

Из теории случайных процессов известно, что реализации такого процесса непрерывны с вероятностью 1 и нигде не дифференцируемы. Следовательно,

данная модель описывает поведение многоэкстремальной непрерывной негладкой функции.

Пусть испытания функции состоят в вычислении ее значений, и проведена серия из k испытаний. Тогда можно вычислить апостериорную (по отношению к поисковой информации $\omega_k = \omega(f, Y_k) = \{(y^i, f^i) : i = 1, \dots, k\}$) плотность распределения $p(z/\omega_k, y)$ неизвестного значения $z=f(y)$. Эта плотность является гауссовой со средним и дисперсией

$$M[z/\omega_k, y] = M[f(y)/\omega_k] = ((y^{i+1} - y)f^i + (y - y^i)f^{i+1}) / (y^{i+1} - y^i) \quad (1.57)$$

$$D[z/\omega_k, y] = D[f(y)/\omega_k] = \sigma(y - y^i)(y^{i+1} - y) / (y^{i+1} - y^i) \quad (1.58)$$

для $y \in [y^i, y^{i+1}]$, где точки измерений считаются упорядоченными по координате в порядке возрастания. Поведение этих характеристик и вид плотности $p(z/\omega_k, y) = p(f(y)/\omega_k)$ представлены на рис.1.22.

В отличие от липшицевой модели (1.52), где для каждого y из области определения D возможные значения $z=f(y)$ лежат в пределах ограниченного промежутка $[f_k^-(y); f_k^+(y)]$, в винеровской модели в точке $y \neq y^i$ ($i=1, \dots, k$) возможны любые значения функции $z=f(y)$, однако более вероятными являются значения мало отличающиеся от математического ожидания $M[f(y)/\omega_k]$. В рамках данной модели, с учетом накопленной поисковой информации, можно оценить, например, вероятность того, что при измерении в точке y будет получено значение функции $f(y) < f_k^* - \varepsilon$, улучшающее рекордное значение f_k^* , достигнутое к k -му шагу, более чем на ε . Это позволяет выделить области

$$D_k = \{y \in D : P(f(y) < f_k^* - \varepsilon / \omega_k) > \delta\}, \quad (1.59)$$

где вероятность вычисления значения функции, более чем на ε улучшающего достигнутое наименьшее ее значение, превосходит заданное $\delta > 0$.

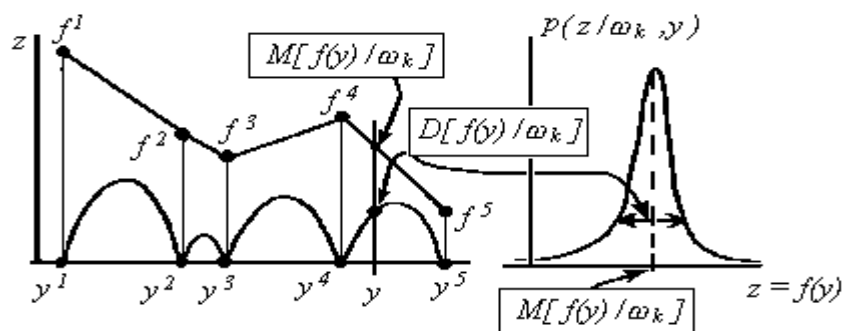


Рис.1.22. Поведение характеристик винеровской модели функции

Множество (1.59) является аналогом множества D_k из (1.56) для липшицевой модели функции. Именно в этом множестве следует размещать очередное испытание задачи.

1.4.2.3. Неполные адаптивные вероятностные модели

В предыдущем пункте был описан классический пример построения вероятностной модели для класса одномерных многоэкстремальных функций. Ясно, что вероятностные модели, по сравнению с детерминированными, определяющими лишь принадлежность функции к некоторому классу $f \in \Phi$, позволяют более гибко и полно учитывать имеющуюся о задаче априорную

информацию и поэтому более привлекательны для построения алгоритмов оптимизации. Однако имеются серьезные «технические» трудности, связанные с использованием этого подхода.

Основных трудностей две: во-первых, сложность построения полного вероятностного описания случайного процесса или поля с нужными свойствами (т.е. согласованного с имеющимися априорными представлениями о задаче), во-вторых, сложность учета полученной поисковой информации (т.е. результатов проведенных испытаний задачи), поскольку этот учет требует пересчета априорных вероятностных распределений в апостериорные по правилам Байеса.

Возможный выход был указан Ю.И. Неймарком [34]. Он предложил использовать так называемые адаптивные стохастические модели, построенные в стиле, характерном для оценок плотностей вероятности в непараметрической статистике. Действительно, при использовании вероятностной модели для построения правила выбора следующей точки испытания y^k достаточно знать только плотность условного вероятностного распределения $p(z/\omega_k, y)$, описывающего вероятности принадлежности неизвестного значения $z=f(y)$ любому заданному измеримому подмножеству числовой оси (в частности, видом этой плотности определяется множество D_k в (1.59)). Таким образом, исходные представления о свойствах функции, в конце концов, преобразуются в определенный вид зависимости этих плотностей вероятности от координат точки y и результатов выполненных испытаний ω_k .

Суть подхода, предложенного Неймарком Ю.И., состояла в том, чтобы ввести прямые правила построения плотностей $p(z/\omega_k, y)$, отказавшись от использования полного вероятностного описания, т.е. задавать только эти плотности. Поскольку пересчет такого неполного вероятностного описания по формулам Байеса невозможен, вводятся правила $\Xi_k(z, \omega_k, y)$ непосредственного построения требуемых распределений вероятности, так что принимается

$$p(z/\omega_k, y) = \Xi_k(z, \omega_k, y), \quad (1.60)$$

где вид зависимости $\Xi_k(z, \omega_k, y)$ постулируется. Одним из естественных требований, налагаемых на получаемые распределения, состоит в том, чтобы для любой точки проведенного испытания y^i ($0 \leq i \leq k$) при $y \rightarrow y^i$ плотность $p(z/\omega_k, y)$ стремилась к δ -функции Дирака $\delta(z - f^i)$, и при удалении y от множества точек испытаний возрастала дисперсия случайной величины, распределенной с плотностью $p(z/\omega_k, y)$. Очевидно, что в форме (1.60) можно представить распределения любого вида, например, подобные представленным на рис.1.22.

Поскольку пересчет по Байесу не используется, фактически речь идет об использовании после каждого испытания новой неполной вероятностной модели, порожденной правилами (1.60). Поэтому *задание адаптивной стохастической (вероятностной) модели состоит в задании системы этих правил.*

Приведем примеры построения адаптивных вероятностных моделей для нескольких классов функций.

Адаптивная вероятностная модель липшицевых функций

Пусть $f \in \Phi = Lip(D)$ с константой L , тогда после накопления поисковой информации ω_k , включающей k измерений функции $f \forall y \in D$ значение функции $z=f(y)$ принадлежит интервалу $[f_k^-(y); f_k^+(y)]$, границы которого определяются в (1.53), (1.54). Будем трактовать (см. [16, 31]) неизвестное значение $z=f(y)$ как

реализацию случайной величины $\xi(y; \omega_k)$, имеющей в этом интервале равномерное распределение. Таким образом модель (1.60) примет вид

$$p(z / \omega_k, y) = \begin{cases} 0, & z \notin [f_k^-(y); f_k^+(y)] \\ 1/(f_k^+(y) - f_k^-(y)), & z \in [f_k^-(y); f_k^+(y)] \end{cases} \quad (1.61)$$

Вид функции математического ожидания $M[f(y)/\omega_k]$, соответствующей этой модели, приведен на рис.1.23.

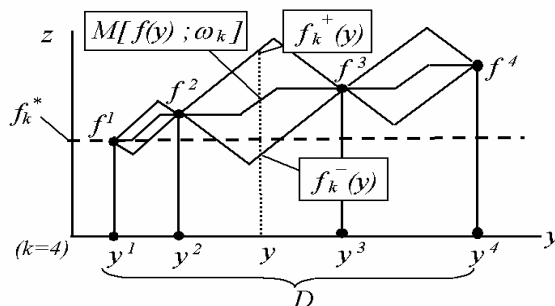


Рис.1.23. Математическое ожидание и границы значений в вероятностной модели липшицевой функции

Преимущество построенной вероятностной модели, по сравнению с детерминированной моделью (1.53)–(1.54), заключается в том, что соответствующие ей оценки подобластей D_k возможного положения решения вида (1.59) используют, в отличие от оценок (1.56), не только нижние, но и верхние границы возможных значений функции $f(y)$. Таким образом, учитываются величины возможных отклонений значений функции от нижних оценок. Модели вида (1.61) будут использованы в главе 6.

Адаптивная вероятностная модель для непрерывных и кусочно-непрерывных функций

Поскольку для функций указанного класса при произвольно выбранном y из D диапазон возможных значений $z=f(y)$ не ограничен, то детерминированные модели вообще не применимы. Для этого класса можно построить адаптивную вероятностную модель. Ее удобнее записать не через плотности $p(z/\omega_k, y)$, а через соответствующие им функции распределения

$$F^f(z / \omega_k, y) = \int_{-\infty}^z p(\tilde{z} / \omega_k, y) d\tilde{z} .$$

Пусть $u(z)$ — функция распределения некоторой случайной величины ξ с нулевым средним и единичной дисперсией, причем $0 < F(z) < 1$. Согласно [31] примем

$$F^f(z / \omega_k, y) = F\left((z - M^f(\omega_k, y)) / D^f(\omega_k, y) \right), \quad (1.62)$$

где

$$M^f(\omega_k, y) = \left(\sum_{i \in I(y)} (f^i + \tilde{L}^f \|y - y^i\|) / \|y - y^i\| \right) / \left(\sum_{i \in I(y)} 1 / \|y - y^i\| \right), \quad (1.63)$$

$$D^f(\omega_k, y) = \sigma_k^f \min \{ \|y - y^i\|^S : i = 1, \dots, k \}. \quad (1.64)$$

В этих соотношениях множество $I(y)$ по разному трактуется для размерности пространства поиска $N > 1$ и $N = 1$. При $N = 1$ под $I(y)$ понимается множество номеров точек испытаний, образующих наименьший интервал, содержащий y , а

при $N > I$ — множество номеров r штук ближайших к y точек y^i проведенных испытаний.

Значение \tilde{L}^f вычисляется по результатам первых \bar{k} испытаний как среднее по всем y^j ($j=1, \dots, \bar{k}$) среди значений

$$\tilde{l}_j = \frac{2}{r(r-1)} \sum_{\substack{i_1 \in I(y^j), i_2 \in I(y^j), \\ i_1 < i_2}} |f^{i_1} - f^{i_2}| / \|y^{i_1} - y^{i_2}\|, \quad (1.65)$$

а параметр σ_k^f оценивается по информации ω_k . Использование ограниченного числа испытаний \bar{k} при вычислении коэффициента \tilde{L}^f предотвращает возможность его неограниченного роста при увеличении k .

Заметим, что математическое ожидание и дисперсия случайной величины $\xi(y)$, соответствующей функции распределения (1.62), будут совпадать со значениями выражений $M^f(\omega_k, y)$ и $D^f(\omega_k, y)$. Их вид иллюстрируется на рис.1.24.

Замечание. Если функция f является кусочно непрерывно-дифференцируемой и при ее испытании в точках y^i кроме значений f^i измеряется еще и ее градиент ∇f^i , то модель (1.62)–(1.64) позволяет легко учесть результаты таких более сложных испытаний. Для этого достаточно изменить вид числителя в выражении (1.63) так, чтобы получилось соотношение, приведенное ниже.

$$M^f(\omega_k, y) = \left(\sum_{i \in I(y)} (f^i + (\nabla f^i, y - y^i)) / \|y - y^i\| \right) / \left(\sum_{i \in I(y)} 1 / \|y - y^i\| \right).$$

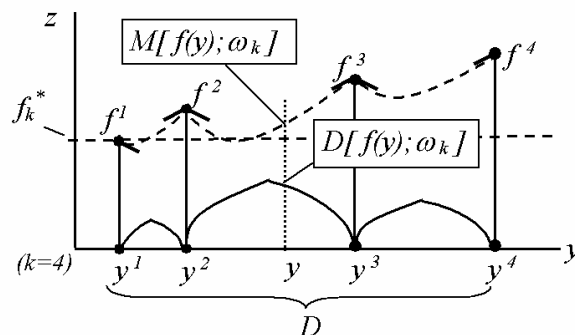


Рис.1.24. Математическое ожидание и дисперсия модели кусочно-непрерывной функции

Адаптивная вероятностная модель функции обобщенного ограничения

Рассмотрим еще один случай, показывающий возможность учета (с помощью адаптивной вероятностной модели) информации о специальной структуре функции. Сделаем это на примере задачи общего вида (1.21)–(1.23), при построении модели которой необходимо учесть набор функций ограничений–неравенств, определенных на множестве $D \subseteq R^N$: $g(y) = (g_1(y), \dots, g_m(y))$,

$$g_1(y) \leq 0, \dots, g_m(y) \leq 0.$$

Будем считать, что $g(y)$ непрерывна или кусочно-непрерывна на D .

Пусть поисковая информация ω_k включает результаты измерений всех компонент функции $g(y)$, т.е. $\omega_k = \omega_k(g) = \omega(g, Y_k)$. Введем обобщенную функцию ограничения

$$G(y) = \max \{g_1(y); \dots; g_m(y)\}. \quad (1.66)$$

Очевидно, что исходная совокупность неравенств эквивалентна одному: $G(y) \leq 0$. При построении некоторых из методов многоэкстремальной оптимизации, рассматриваемых в главе 4, понадобится вероятностная модель обобщенной функции ограничения, позволяющая учесть специфику ее структуры, определяемую видом соотношения (1.66). На первый взгляд может показаться, что можно поступить достаточно просто, пересчитав поисковую информацию $\omega_k(g) = \omega(g, Y_k)$ в поисковую информацию $\omega_k(G) = \omega(G, Y_k)$ для функции $G(y)$ и воспользовавшись вероятностной моделью (1.62)–(1.64), построив ее по $\omega_k = \omega_k(G) = \omega(G, Y_k)$. Однако такой подход приведет к частичной потере информации о структуре функции G из (1.66).

Укажем другой способ получения модели, учитывающий структуру G . Вначале построим аппроксимации отдельно для каждой из функций $g_j(y)$ по результатам их измерений $\omega_k = \omega_k(g) = \omega(g, Y_k)$:

$$M^{g_j}(\omega_k(g), y) = \left(\sum_{i \in I(y)} g_j^i / \|y - y^i\| \right) / \left(\sum_{i \in I(y)} 1 / \|y - y^i\| \right). \quad (1.67)$$

Далее, используя функцию распределения $u(z)$ с нулевым средним и единичной дисперсией, такую что $0 < F(z) < 1$, построим адаптивную вероятностную модель $G(y)$ с функцией распределения вида

$$F^G(z / \omega_k(g), y) = F\left((z - M^G(\omega_k(g), y)) / D^G(\omega_k(g), y) \right), \quad (1.68)$$

$$M^G(\omega_k(g), y) = \max \{ M^{g_1}(\omega_k(g), y); \dots; M^{g_m}(\omega_k(g), y) \}, \quad (1.69)$$

$$D^G(\omega_k(g), y) = \sigma_k^G \min \{ \|y - y^i\|^{\gamma} : i = 1, \dots, k \}. \quad (1.70)$$

На рис.1.25 показано поведение математического ожидания $M[G(y); \omega_k(g)] = M^G(\omega_k(g), y)$ и дисперсии обобщенного ограничения G при конкретных результатах измерений двух функций ограничений g_1 и g_2 .

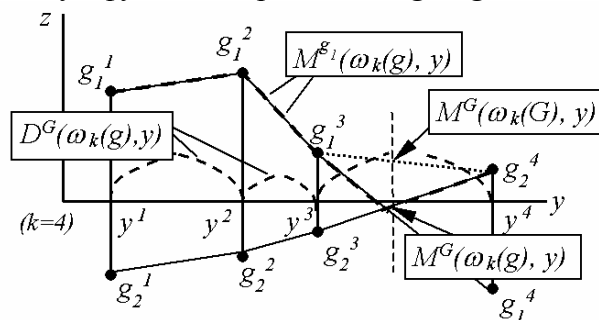


Рис.1.25. Математическое ожидание и дисперсия модели функции обобщенного ограничения

Построенная модель обобщенной функции ограничения, как уже отмечалось выше, будет использована в главе 4.

👉 Контрольные вопросы и упражнения

1. Укажите множество значений y на рис.1.25, в которых вероятность выполнения ограничения $G(y) \leq 0$ равна нулю.
2. Считая, что в (1.68) в качестве $F(z)$ используется функция нормального распределения, укажите на рис.1.25 подынтервалы значений y , для которых вероятность выполнения ограничения $G(y) \leq 0$ не менее 0,5.
3. Как изменится результат в упражнении 2, если в модели (1.68) вместо

$M^G(\omega_k(g), y)$ из (1.69) использовать выражение другого вида, не учитывающее особенность структуры (1.66), а именно — из (1.71) (см. рис.1.25).

$$M^G(\omega_k(G), y) = \left(\sum_{i \in I(y)} G^i / \|y - y^i\| \right) / \left(\sum_{i \in I(y)} 1 / \|y - y^i\| \right). \quad (1.71)$$

Заметим, что к настоящему времени разработано несколько алгоритмов, основанных на использовании адаптивных вероятностных моделей рассмотренных типов, начиная от простых задач и до более сложных (см., например, [16]).

Лист регистрации изменений

Дата	Автор	Комментарии
19.05.02	Городецкий С.Ю.	Создание документа, набор р. 1.4
21.05.02	Городецкий С.Ю.	Набор разделов 1.1-1.2
10.07.02	Городецкий С.Ю.	Внесение изменений в 1.1-1.2
14.07.02	Городецкий С.Ю.	Внесение изменений в 1.3
22.12.02	Городецкий С.Ю.	Окончательная редакция 1.1-1.3
28.07.02	Городецкий С.Ю.	Внесение изменений в 1.4
11.08.03	Городецкий С.Ю.	Дополнение и переработка 1.4
12.08.03	Городецкий С.Ю.	Добавление рисунков 1.4
04.09.03	Городецкий С.Ю.	Корректурa + новый формат
09.09.03	Городецкий С.Ю.	Исправление ссылок

Глава 2. Теоретические основы аналитического решения задач оптимизации

В этой главе будут изучены две темы, логически связанные между собой. Первая тема посвящается выводу условий экстремума в достаточно общей форме, так чтобы они были применимы к многокритериальным задачам. Здесь приводятся естественные обобщения классических результатов по условиям экстремума для однокритериальных задач математического программирования (функция Лагранжа, условия Куна–Таккера и близкие к ним) на многокритериальные задачи.

Эти условия можно использовать как для аналитического отыскания решений в конкретных задачах, имеющих относительно простую структуру, так и для численно–аналитического решения более сложных задач с применением математических пакетов таких, например, как Mathcad.

Следует отметить, что знание теории по условиям экстремума важно для понимания теории двойственности, лежащей в основе ряда вычислительных методов, используемых при решении экстремальных задач.

Материал по условиям экстремума приводится с полным теоретическим обоснованием, поскольку для его глубокого понимания необходимо не только знание самих условий, но и способов их получения.

Вторая тема, рассматриваемая в главе, — элементы теории двойственности, основанной на использовании функции Лагранжа при построении двойственной задачи. Этот материал тесно связан с вычислительными методами, основанными на модифицированных функциях Лагранжа, для решения задач с ограничениями.

2.1. Обобщение условий экстремума на задачи векторной оптимизации

В предыдущей главе задача векторной оптимизации была поставлена в следующей форме (будем кратко называть ее задачей FYD):

- Задача FYD

$$f(y) \rightarrow \min, y \in Y \subseteq R^N, f: Y \rightarrow R^n, \quad (2.1)$$

$$Y = \{y \in D : g(y) \leq 0, h(y) = 0\}, \quad (2.2)$$

$$D = \{y \in R^N : a \leq y \leq b\}, g: D \rightarrow R^m, h: D \rightarrow R^p \quad (2.3)$$

Все неравенства, записанные в векторной форме, понимаются покомпонентно. Компоненты векторов нумеруются нижним индексом, например, $y = (y_1, \dots, y_N)$, $f = (f_1, \dots, f_n)$ и т.д.

Форма записи (2.1)–(2.3) удобна при построении вычислительных методов, но не удобна при теоретическом анализе. Целесообразно перейти к иной форме представления этой задачи, заменив в (2.2) требование $y \in D$ на $y \in E$, считая, что новое множество E может включать в себя часть ограничений в форме равенств и неравенств из (2.2). В то же время, часть двусторонних ограничений на переменные из (2.3) может быть переведена в форму неравенств общего вида. Задача получит следующее представление.

а) Задача FVE

$$f(y) \rightarrow \min, y \in Y, f: Y \rightarrow R^n, \quad (2.4)$$

$$Y = \{y \in E \subseteq R^N : g(y) \leq 0, h(y) = 0\}, \quad (2.5)$$

$$g: E \rightarrow R^m, h: E \rightarrow R^p$$

Подчеркнем, что вектор–функции ограничений g, h и их размерности m и p в (2.5) могут отличаться от тех, которые представлены в исходной постановке (2.2), (2.3).


Таким образом, постановка (2.4), (2.5) дает более гибкую форму представления исходной задачи, позволяя вынести на уровень явной формы записи любые ограничения (включая координатные), скрыв в множестве E все остальные. Если, например, в множество E «спрятать» все ограничения–равенства, то в (2.5) они явно вообще не будут присутствовать.

Перейдем к получению условий экстремума. Построим аналог функции Лагранжа для задач векторной оптимизации. Как было показано в методе главного критерия (пункт 1.3.3 первой главы), для вычисления конкретных слабо эффективных решений все компоненты целевой вектор–функции, кроме одного, можно перевести в разряд ограничений–неравенств. Это наводит на мысль, что компоненты функции $f(y)$ должны входить в функцию Лагранжа в той же форме, что и компоненты $g(y)$.

Определение 2.1. Функцию вида

$$L(y, v, \lambda, \mu) = (f(y), v) + (g(y), \lambda) + (h(y), \mu) \quad (2.6)$$

Назовем функцией Лагранжа для задачи (2.4), (2.5).

 **Замечание.** Основная идея, связанная с введением функции Лагранжа, состоит в том, чтобы свести проблему решения (многокритериальной) задачи с ограничениями к задаче (однокритериальной) без таких ограничений.

Будем искать условия слабой эффективности точки $y^0 \in Y$.

Аналог условий Куна-Таккера¹ имеет вид

$$y \in Y - \text{условие допустимости}, \quad (2.7)$$

$$\exists (v^0, \lambda^0, \mu^0) \neq 0 - \text{условие нетривиальности}, \quad (2.8)$$

$$v^0 \geq 0, \lambda^0 \geq 0 - \text{условие неотрицательности}. \quad (2.9)$$

Принцип минимума:

$$L(y^0, v^0, \lambda^0, \mu^0) = \min \{L(y, v^0, \lambda^0, \mu^0) : y \in E\}, \quad (2.10)$$

$$\lambda^0_i g_i(y^0) = 0 - \text{условие дополняющей нежесткости}. \quad (2.11)$$

Обоснование этих условий будет приведено позднее. Сейчас следует обратить внимание на то, что минимум по множеству (2.5), включающему явно присутствующие ограничения равенства и неравенства, заменяется в (2.10) на

¹ Условия экстремума для многокритериальных задач исследовались в работах многих авторов (см., например, библиографию в [6]). Поэтому использованное название является условным и лишь подчеркивает связь с классическими результатами математического программирования.

операцию взятия минимума только по множеству E . Это происходит за счет того, что явные ограничения включены в функцию Лагранжа (2.6). Кроме того, условие дополняющей нежесткости (2.11) означает, что для допустимых точек y^0 , для которых i -е неравенство выполняется строго (такие неравенства называют *неактивными* в точке y^0), обязательно $\lambda_i^0 = 0$, т.е. такие ограничения в функцию Лагранжа, фактически, не входят.

Заметим, что условия (2.7) – (2.11), вообще говоря, могут выполняться при $\nu^0 = 0$, что означает независимость получаемой в этом случае точки y^0 от целевой вектор–функции. Такие случаи следует считать вырожденными. При этом говорят, что задача *не регулярна* в точке y^0 . В последующем будут получены достаточные условия регулярности задач вида FУЕ.

Теорема 2.1 (достаточное условие слабой эффективности). Если для точки y^0 выполнены условия (2.7) – (2.11) при $\nu^0 \neq 0$, то y^0 – слабо эффективная точка задачи FУЕ ((2.4)-(2.5)).

ДОКАЗАТЕЛЬСТВО выполним от противного. Пусть $\exists y' \in Y$, что $f(y') < f(y^0)$, где y^0 удовлетворяет аналогу условий Куна-Таккера при $\nu^0 \neq 0$. Тогда возможна следующая оценка

$$L(y', \nu^0, \lambda^0, \mu^0) = (f(y'), \nu^0) + (\lambda^0, g(y')) + (\mu^0, h(y')) \leq (f(y'), \nu^0) \leq (f(y^0), \nu^0) = L(y^0, \nu^0, \lambda^0, \mu^0)$$

\leq
 $= 0$
 $\nu^0 \geq 0, \nu^0 \neq 0$

Это противоречит принципу минимума (2.10).

Замечание. Для задач достаточно общего вида условия (2.7)-(2.11) при $\nu^0 \neq 0$ только достаточны, т.е. совсем не обязательно, чтобы они выполнялись для всех слабо эффективных точек.

Можно построить пример, когда для некоторой слабо эффективной точки $y^0 \in Y^0$ условия (2.7)–(2.11) выполняться не будут. Действительно, функция Лагранжа (2.6) совпадает с функцией Лагранжа, соответствующей той же задаче FУЕ, но со скалярной целевой функцией, порожденной линейной сверткой вида $\Psi_\nu(y) = (f(y), \nu)$ (см. пункт 1.3.6 главы 1). А она, как было показано в свойстве 1.4, не позволяет отыскивать все точки из Y^0 , если F — образ множества Y ,

$$F = \{z = f(y) : y \in Y\}$$

не является выпуклым. Обеспечим теперь выпуклость этого множества.

Теорема 2.2 (необходимые условия слабой эффективности). Если E выпукло, f, g – выпуклы на E , а h – аффинная (т.е. $h(y) = Ay + c$) на E , то для того, чтобы y^0 было слабо эффективной в задаче FУЕ необходимо, чтобы выполнялись условия (2.7) – (2.11).

ДОКАЗАТЕЛЬСТВО. Покажем, что обоснование этой теоремы можно изящно свести к линейной отделимости некоторого выпуклого множества S и некоторой его граничной точки P , а сама функция Лагранжа возникает из уравнения разделяющей гиперплоскости.

Введем множества

$$\begin{aligned} S(y) &= \{z = (z_f, z_g, z_h) : z_f \geq f(y), z_g \geq g(y), z_h = h(y)\} \\ S &= \bigcup_{y \in E} S(y) \end{aligned} \tag{2.12}$$

и точку

$$P = \{z = (z_f, z_g, z_h) : z_f = f(y^0), z_g = 0, z_h = 0\}$$

Покажем выпуклость S и то, что $P \in \partial S$ — границе S . Пусть $z^1, z^2 \in S$ и им соответствуют точки $y_1, y_2 \in E$. Обозначим $z^\alpha = (z_f^\alpha, z_g^\alpha, z_h^\alpha) = \alpha z^1 + (1-\alpha)z^2$ и проверим выполнение условия $z^\alpha \in S(y^\alpha)$, где $y^\alpha = \alpha y^1 + (1-\alpha)y^2$. Для этого надо убедиться в выполнении соответствующих неравенств.

Из выпуклости следует, что

$$f(y^\alpha) \leq \alpha f(y^1) + (1-\alpha)f(y^2) \leq z_f^\alpha,$$

аналогично показывается, что $g(y^\alpha) \leq z_g^\alpha$. В силу афинности,

$$h(y^\alpha) = \alpha h(y^1) + (1-\alpha)h(y^2) = \alpha z_h^1 + (1-\alpha)z_h^2 = z_h^\alpha.$$

Таким образом, выпуклость S доказана.

Заметим теперь, что $P \in S(y^0) \subseteq S$. Введем точку $P_\varepsilon = (f(y^0) - \varepsilon \cdot e, 0, 0)$, где

$$e = (1, \dots, 1)^T \in R^n, \varepsilon > 0.$$

Если бы $P_\varepsilon \in S$, то нашлась бы точка $y_\varepsilon \in E$, с

$$f(y_\varepsilon) \leq f(y^0) - \varepsilon \cdot e, \quad g(y_\varepsilon) \leq 0, \quad h(y_\varepsilon) = 0.$$

Но тогда $y_\varepsilon \in Y$ и $f(y_\varepsilon) < f(y^0)$, что невозможно для слабо эффективной точки y^0 . Таким образом, для сколь угодно малого положительного $\varepsilon > 0$ оказывается, что $P_\varepsilon \notin S$, но при $\varepsilon \rightarrow 0$ $P_\varepsilon \rightarrow P \in S$, следовательно P — граничная точка выпуклого S и их можно линейно отделить.

Запишем вектор нормали разделяющей гиперплоскости как $(v^0, \lambda^0, \mu^0) \neq 0$. Тогда $\forall z = (z_f, z_g, z_h) \in S$:

$$(z_f, v^0) + (z_g, \lambda^0) + (z_h, \mu^0) \geq (f(y^0), v^0) + (0, \lambda^0) + (0, \mu^0).$$

Представляя $z_f = f(y) + \varepsilon$, $z_g = g(y) + \delta$, $z_h = h(y)$ ($\varepsilon \geq 0$, $\delta \geq 0$), получим, что $\forall y \in E, \forall \varepsilon \geq 0, \delta \geq 0$

$$(f(y) + \varepsilon, v^0) + (g(y) + \delta, \lambda^0) + (h(y), \mu^0) \geq (f(y^0), v^0) \quad (2.13)$$

Теперь получим из этого неравенства все утверждения, которые нужно обосновывать в теореме.

Выберем $y = y^0$, тогда $\forall \varepsilon \geq 0, \delta \geq 0: (\varepsilon, v^0) + (g(y^0) + \delta, \lambda^0) \geq 0$. Из того, что $g(y^0) \leq 0$ и полученного неравенства следует $v^0 \geq 0, \lambda^0 \geq 0$.

Полагая $y = y^0, \varepsilon = \delta = 0$, из (2.13) получаем $(g(y^0), \lambda^0) \geq 0$. Но $g_i(y^0) \lambda_i \leq 0$ ($i = 1, \dots, n$), откуда получаем условие (2.11).

Выбирая $\varepsilon = 0, \delta = 0$ из (2.13) получаем условие (2.10). Таким образом, теорема доказана.

Сопоставив формулировки теорем 2.1 и 2.2 нетрудно понять, что для превращения обобщенных условий Куна–Таккера (2.7)–(2.11) в критерий слабой эффективности нужно рассмотреть такой подкласс выпуклых задач, где не возникает вырождения условий экстремума, то есть $v^0 \neq 0$.

Наложим ограничения на допустимую область Y , чтобы всегда при любой вектор–функции $f(y)$, имеющей в $y^0 \in Y$ слабо эффективную точку, при записи условий экстремума $f(y)$ на Y для точки y^0 значение v^0 можно было выбирать отличным от 0.

Определение 2.2. Область Y , удовлетворяющая указанным выше условиям для точки y^0 , называется *регулярной в этой точке*. Область Y , регулярная во всех своих точках, называется *регулярной*.

Существует простое достаточное условие регулярности выпуклых областей полученное Слейтером.

Теорема 2.3. (достаточное условие регулярности Слейтера) Пусть E – выпуклое множество, $g(y)$ – выпуклы (покомпонентно) на E , а $h(y)$ – афинны. Тогда если $\exists \bar{y} \in E$, что $h(\bar{y}) = 0$, $g(\bar{y}) < 0$ и кроме того

$$0 \in \text{int } h(E), \quad h(E) = \{z_h = h(y) : y \in E\} \quad (2.14)$$

(где int обозначает внутренность множества), то множество Y будет регулярно во всех своих точках.

Следует обратить внимание на то, что неравенства в точке \bar{y} должны выполняться строго. Заметим также, что при отсутствии ограничений–равенств условие (2.14) не нужно.

Доказательство проведем от противного. Пусть утверждение не верно и существует вектор–функция $f(y)$ со слабо эффективной точкой y^0 , для которой условия экстремума (2.7)–(2.11) выполняются при $v^0 = 0$.

Используя в правой части (2.10) $y = \bar{y}$, получим неравенство

$$(v^0, f(y^0)) + (g(y^0), \lambda^0) + (\mu^0, h(y^0)) \leq (v^0, f(\bar{y})) + (\lambda^0, g(\bar{y})) + (\mu^0, h(\bar{y})),$$

$\boxed{v^0=0}$ $\boxed{=0 \text{ по (2.11)}}$ $\boxed{=0}$

или

$$0 \leq (\lambda^0, g(\bar{y})) + (\mu^0, h(\bar{y})), \text{ где } (\lambda^0; \mu^0) \neq 0.$$

$\boxed{<0}$ $\boxed{=0}$

Если $\exists \lambda_i^0 > 0$, то получаем противоречие. Если $\lambda^0 = 0$, то $\mu^0 \neq 0$. В этом случае, за счет выбора $y' \in O_\varepsilon(\bar{y}) \cap E$, обеспечим знак компонент вектора $h(y)$, противоположный знакам компонент в векторе μ^0 . Это возможно, т.к. $h(\bar{y}) = 0 \in \text{int } h(E)$. При этом мы вновь придем к противоречию со знаком неравенства.

Следствие 2.3.1. При выполнении условий регулярности Слейтера для допустимой области и выпуклости компонент вектор–функции f на выпуклом E обобщенные условия Куна–Таккера (2.7)–(2.11) являются необходимыми и достаточными условиями слабо эффективной точки y^0 .

Существует несколько форм записи условий (2.7)–(2.11). Одна из них – в терминах седловой точки функции Лагранжа. Позднее мы будем обращаться к этой форме записи при построении вычислительных методов учета ограничений, основанных на использовании модифицированных функций Лагранжа.

Теорема 2.4. (о записи условий оптимальности через седловую точку функции Лагранжа) При $v^0 \neq 0$ условия (2.7)-(2.11) выполняются для точки y^0 тогда и только тогда, когда

$$y^0 \in Y, \exists (v^0, \lambda^0, \mu^0) \neq 0, v^0 \geq 0, \lambda^0 \geq 0 \quad (2.15)$$


что $\forall y \in E \quad \forall \lambda \geq 0, \mu$

$$L(y^0, v^0, \lambda, \mu) \leq L(y^0, v^0, \lambda^0, \mu^0) \leq L(y, v^0, \lambda^0, \mu^0) \quad (2.16)$$

ДОКАЗАТЕЛЬСТВО. Легко видеть, что правое неравенство в (2.16) совпадает с (2.10), а (2.15) – с (2.8)-(2.9). Остается показать эквивалентность левого неравенства в (2.16) и условия дополняющей нежесткости (2.11).

С учетом допустимости $y^0 \in Y, h(y^0) = 0$, поэтому левое неравенство (2.16) примет вид: $\forall \lambda \geq 0 \quad (\lambda, g(y^0)) \leq (\lambda^0, g(y^0))$.

Поскольку $g(y^0) \leq 0$, оно эквивалентно (2.11).

 **Замечание.** Условие (2.16) позволяет при фиксированном $v^0 \geq 0, v^0 \neq 0$ свести поиск соответствующей слабо эффективной точки y^0 в задаче с функциональными ограничениями (2.4)-(2.5) к отысканию седловой точки функции Лагранжа в области $y \in E, \lambda \geq 0$ пространства переменных (y, λ, μ) .

Этот факт используется в некоторых общих вычислительных методах учета ограничений.

2.2. Условия оптимальности в дифференциальной форме для многокритериальных задач оптимизации специального и общего вида

Для аналитического и численно-аналитического решения задач оптимального выбора более удобной является дифференциальная форма записи условий (2.7)-(2.11).

2.2.1. Условия первого порядка

Пусть функции f, g – дифференцируемы и выпуклы на выпуклом E, h – аффинная, а область Y регулярна, тогда для $y^0 \in Y^0$ (множество слабо эффективных точек) при дополнительном предположении о том, что y^0 – внутренняя для E ($y^0 \in \text{int } E$) из теоремы 2.2 следует, что найдется $(v^0, \lambda^0, \mu^0) \neq 0, v^0 \geq 0, \lambda^0 \geq 0, \lambda_i^0 g_i(y^0) = 0$ ($i=1, \dots, m$), что

$$\nabla_y L(y^0, v^0, \lambda^0, \mu^0) = 0. \quad (2.17)$$

Покажем, что эти условия будут достаточными для слабой эффективности y^0 . Действительно, функция $L(y, v^0, \lambda^0, \mu^0)$ выпукла на выпуклом E и дифференцируема, а ее градиент по y в точке y^0 обращается в 0. Тогда по критерию выпуклости дифференцируемой функции

$\forall y \in E \quad L(y, v^0, \lambda^0, \mu^0) \geq L(y^0, v^0, \lambda^0, \mu^0) + (0, y - y^0)$ в силу (2.17), следовательно, y^0 удовлетворяет принципу минимума (2.10). Из теоремы (2.1) следует слабая эффективность точки y^0 . Таким образом ДОКАЗАНА следующая теорема.

Теорема 2.5. Если множество E выпукло, f, g дифференцируемы на E и выпуклы, а h – афинна, тогда при регулярности допустимой области Y для слабой эффективности точки y^0 , являющейся внутренней точкой множества E , необходимо и достаточно, чтобы

$$y^0 \in Y, \exists (v^0, \lambda^0, \mu^0), v^0 \geq 0, v^0 \neq 0, \lambda^0 \geq 0, \quad (2.18)$$

что

$$(\nabla f(y^0), v^0) + \sum_{j \in J(y^0)} \lambda_j^0 \nabla g_j(y^0) + (\mu^0, \nabla h(y^0)) = 0 \quad (2.19)$$

Через $J(y^0)$ обозначено множество номеров активных ограничений–неравенств:

$$J(y^0) = \{j \in \{1, \dots, n\} : g_j(y^0) = 0\}.$$

Условия теоремы 2.5 весьма удобны для аналитического отыскания слабо эффективных точек, либо численно–аналитического их определения с помощью пакетов для математических расчетов. Удобство состоит в том, что вопрос об отыскании точки y^0 сводится к решению систем уравнений вида (2.19) в предположении, что $J(y^0) = J$ (J — рабочий набор активных ограничений–неравенств), с добавленными к ним условиями

$$g_{j_i}(y^0) = 0 \quad (i = 1, \dots, r), \quad \{j_1, \dots, j_r\} = J \quad (2.20)$$

и уравнениями ограничений–равенств

$$h_s(y^0) = 0 \quad (s = 1, \dots, p), \quad (2.21)$$

которые вместе с (2.20) дают $r+p$ дополнительных уравнений.


В результате для каждого значения $v^0 \geq 0$, возникает система из $N+r+p$ уравнений (в общем случае, нелинейных) относительно $N+r+p$ неизвестных $y_1^0, \dots, y_N^0, \lambda_{j_1}^0, \dots, \lambda_{j_r}^0, \mu_1^0, \dots, \mu_p^0$.

Основная сложность в их использовании заключается в том, что кроме уравнений системы (2.19)–(2.21) решение должно удовлетворять набору неравенств

$$g_{j_s}(y^0) \leq 0 \quad \forall j_s \notin J \quad (2.22)$$

$$\lambda_{j_i} \geq 0 \quad \forall j_i \in J \quad (2.23)$$

В случае нарушения одного из них необходимо изменить множество J и повторить все действия еще раз. Естественно, что при коррекции рабочего множества J (номеров активных ограничений) следует использовать информацию о нарушении условий (2.22), (2.23). Если точка y^0 не удовлетворяет некоторым из неравенств (2.22), часть из них следует добавить в рабочий набор J , а если оказалось $\lambda_{j_i}^0 < 0$, то соответствующий номер j_i следует вывести из J .

 **Замечание.** Последнее правило означает, что после проверки гипотезы $J(y^0) = J$, следует сместить точку y^0 с границ тех активных ограничений, где $\lambda_{j_i}^0 < 0$ в область с $g_{j_i}(y) < 0$. Это приведет к уменьшению значений функции линейной свертки $\Psi_{v^0}(y)$ из (1.29) векторного критерия $f(y)$.

ДОКАЗАТЕЛЬСТВО. Это правило основывается на следующем наблюдении. Построим вектор w смещения из точки y^0 так, чтобы $(\nabla g_{j_i}(y^0), w) < 0$ для $\lambda_{j_i}^0 < 0$, $(\nabla g_{j_i}(y^0), w) = 0$ для $\lambda_{j_i}^0 \geq 0$ и $(\nabla h_s(y^0), w) = 0$. По построению вектор w направлен внутрь допустимой области по отношению к ограничениям–неравенствам с отрицательными значениями множителей Лагранжа. Будем смещать точку y так, чтобы отклонение $y - y^0$ зависело от параметра сдвига α как $y^\alpha - y^0 = w \cdot \alpha + o(\alpha)$, а поправка $o(\alpha)$ не уводила бы не только с аффинных ограничений–равенств (т.е. $h(y^\alpha) \equiv 0$), но и с границ тех ограничений–неравенств, которые остаются активными, т.е. $g_{j_i}(y^\alpha) \equiv 0$ для $\lambda_{j_i}^0 \geq 0$.


Умножая скалярно на вектор w основное равенство (2.19) из последней теоремы, получим, что

$$\left(\sum_{i=1}^n \nabla f_i(y^0), v_i^0, w \right) = - \sum_{\lambda_{j_i}^0 < 0} (\nabla g_{j_i}(y^0), w) \lambda_{j_i}^0 < 0.$$

Если рассмотреть функцию $\Psi_{v^0}(y)$ линейной свертки векторного критерия f (см. пункт 1.3.6 первой главы), то, используя в качестве весовых коэффициентов v_1^0, \dots, v_n^0 , последнее неравенство можно переписать в новой форме, используя производную по направлению w : $\frac{\partial \Psi_{v^0}(y^0)}{\partial w} < 0$. Таким образом, линейная свертка $\Psi_{v^0}(y)$ строго убывает при малом смещении вдоль w от точки y^0 , а поскольку $\Psi_{v^0}(y^\alpha) - \Psi_{v^0}(y^0) = \frac{\partial \Psi_{v^0}(y^0)}{\partial w} \cdot \alpha + o(\alpha)$, то это свойство сохранится при малых значениях α и для смещений по кривой y^α .

Поскольку из материала раздела 1.3 следует, что слабо эффективная точка, соответствующая выбранным значениям множителей v^0 , всегда может быть найдена в выпуклой задаче минимизацией функции линейной свертки $\Psi_{v^0}(y)$, то смещение вдоль y^α будет способствовать, в указанном смысле, улучшению построенной оценки решения. Таким образом, утверждение верно.

Следует обратить внимание на то, что смещение вдоль y^α не обязано порождать точки с $f(y^\alpha) \leq f(y^0)$. Часть критериев при малых смещениях будет строго убывать, но другая часть может возрастать (см. рис. 1.12).

 **Замечание.** Следует обратить внимание на то, что достаточность условий (2.18), (2.19) может быть обоснована при гораздо более слабых ограничениях, чем это сделано в теореме 2.5.

Приведем эти условия в виде отдельной теоремы.

Теорема 2.6. Пусть E – выпукло, f – псевдовыпукла, g – дифференцируемы на E , $y^0 \in \text{int } E$ и $\forall j \in J(y^0)$: $g_j(y)$ – квазिवыпуклы, а $h_s(y)$ ($s=1, \dots, p$) – квазिवыпуклы и квазивогнуты, тогда выполнение условий (2.18), (2.19) достаточно для слабой эффективности y^0 .

ДОКАЗАТЕЛЬСТВО. Пусть это не так и $\exists y^1 \in Y$, что $f(y^1) < f(y^0)$. Поскольку $\forall j \in J(y^0)$ $g_j(y^1) \leq 0 = g_j(y^0)$, то из свойства 1.7'' квазिवыпуклой функции

(см. пункт 1.4.1) следует, что $(\nabla g_j(y^0), y^1 - y^0) \leq 0$. Аналогично, из свойств функции $h(y)$ получим $(\nabla h(y^0), y^1 - y^0) = 0$.

Домножая скалярно равенство (2.19) на вектор направления $d = y^1 - y^0$, получим

$$\sum_{i=1}^n v_i^0 (\nabla f_i(y^0), d) \geq 0, \text{ что при } v^0 \geq 0 \text{ и } v^0 \neq 0 \text{ означает, что } \exists i : (\nabla f_i(y^0), d) \geq 0. \text{ Из}$$

псевдовыпуклости $f_i(y^0)$ и последнего неравенства следует $f_i(y^1) \geq f_i(y^0)$. Но это противоречит нашему предположению.

Таким образом, использование введенных в разделе 1.4 моделей поведения функций позволяет расширить область применения теоремы 2.5.

Дальнейшее ее обобщение требует отказа от любых форм выпуклости. При этом полученные условия, очевидно, не смогут обеспечить слабой эффективности точки y^0 , и только ее локальную слабую эффективность. Это следует из того, что выполнение условий вида (2.19) отражает только локальные свойства функций задачи, а слабая эффективность точки является ее интегральной характеристикой в области Y .

Представляет интерес способ обоснования приведенной ниже теоремы. По характеру он является более наглядным и геометрическим, чем в теореме 2.2. А именно, условия экстремума вначале формулируются в терминах непересекаемости некоторых множеств, а затем сводятся к непересекаемости их конических аппроксимаций, и только потом условие непересекаемости приводится к аналитической форме записи в градиентной форме.

Теорема 2.7. Пусть $f(y)$, $g(y)$ -дифференцируемы в точке y^0 , а $h(y)$ -непрерывно дифференцируема в y^0 . Чтобы точка y^0 , принадлежащая $\text{int } E$, была локально слабо эффективной (см. определение 1.5 из раздела 1.3) необходимо, чтобы

$$\exists (v^0, \lambda^0, \mu^0) \neq 0, \text{ что } v^0 \geq 0, \lambda^0 \geq 0, \quad (2.24)$$

$$\nabla_y L(y^0, v^0, \lambda^0, \mu^0) = 0, \quad (2.25)$$

$$\forall i = 1, \dots, m \quad \lambda_i^0 \cdot g_i(y^0) = 0. \quad (2.26)$$

ДОКАЗАТЕЛЬСТВО. Исключим из рассмотрения тривиальный случай, когда $\nabla h_1(y^0), \dots, \nabla h_p(y^0)$ –линейно зависимая система. В этом случае утверждение теоремы будет элементарно выполнено при $v^0 = 0, \lambda^0 = 0, \mu^0 \neq 0$.

Далее будем считать $\nabla h_1(y^0), \dots, \nabla h_p(y^0)$ линейно независимыми. Рассмотрим следующие множества

$$H = \{y \in R^N : h(y) = 0\}, \quad G_j = \{y \in R^N : g_j(y) \leq 0\} \quad (j = 1, \dots, m)$$

$$F_i(y^0) = \{y \in R^N : f_i(y) < f_i(y^0)\} \quad (i = 1, \dots, n)$$

Условия экстремума с использованием этих множеств можно записать в следующей очевидной форме: для того, чтобы точка y^0 , внутренняя для E , была локально слабо эффективной, необходимо существование $\varepsilon > 0$ при котором

$$\bigcap_{i=1}^n F_i(y^0) \cap \bigcap_{j=1}^m G_j \cap H \cap O_\varepsilon(y^0) = \emptyset. \quad (2.27)$$

Введем в окрестности точки y^0 аппроксимации множеств, входящих в условие (2.27), конусами:

$$K^H(y^0) = \{d \in R^N : (d, \nabla h_s(y^0)) = 0, s = 1, \dots, p\}, \quad (2.28)$$

$$K_j^g(y^0) = \{d \in R^N : (d, \nabla g_j(y^0)) < 0\}, j \in J(y^0), \quad (2.29)$$

$$K_i^f(y^0) = \{d \in R^N : (d, \nabla f_i(y^0)) < 0\}, (i = 1, \dots, n). \quad (2.30)$$

Лемма 2.1. Если выполнено условие непересекаемости (2.27), то

$$\underbrace{\bigcap_{i=1}^n K_i^f(y^0)}_{K^f} \cap \underbrace{\bigcap_{j \in J(y^0)} K_j^g(y^0)}_{K^Y} \cap K^H(y^0) = \emptyset. \quad (2.31)$$

Рис.2.1 иллюстрирует утверждение леммы.

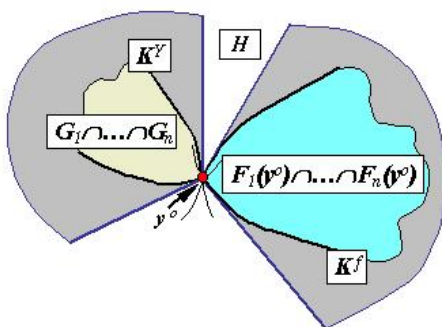


Рис.2. 1. Соотношение между коническими аппроксимациями множеств

ДОКАЗАТЕЛЬСТВО леммы. Пусть $\exists d \in K^f \cap K^Y$. Тогда $\forall i=1, \dots, n: (\nabla f_i(y^0), d) < 0$, $\forall j \in J(y^0): (\nabla g_j(y^0), d) < 0$ и $d \in K^H$. В условиях теоремы найдется $o(\varepsilon)$, что $h(y^0 + \varepsilon d + o(\varepsilon)) = 0$ (рис.2.2).

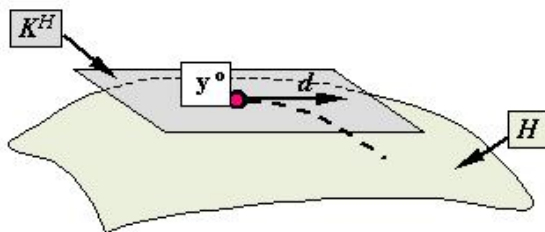


Рис.2.2. Кривая y^ε , порожденная на многообразии H касательным направлением d

Обозначим $y^\varepsilon = y^0 + \varepsilon d + o(\varepsilon)$. Тогда при малых ε

$$f(y^\varepsilon) = f(y^0) + (\nabla f(y^0), d) \varepsilon + o(\varepsilon) < f(y^0).$$

Аналогичные рассуждения показывают, что при достаточно малых $\varepsilon: g_j(y^\varepsilon) < 0$, $j \in J(y^0)$, и $g_s(y^\varepsilon) < 0$ при $s \notin J(y^0)$ за счет того, что $g_s(y^0) < 0$.

В результате оказалась найдена точка $y^\varepsilon \in O_\varepsilon(y^0) \cap Y$ со значением $f(y^\varepsilon) < f(y^0)$, что противоречит слабой эффективности точки y^0 . Лемма доказана.

Нам осталось показать, что условие (2.31) равносильно (2.24)–(2.26). Это можно сделать, сконструировав специальные вспомогательные линейные задачи оптимизации. Поскольку мы рассматриваем случай линейной независимости системы векторов $\nabla h_1(y^0), \nabla h_2(y^0), \dots, \nabla h_p(y^0)$, то конус $K^H(y^0)$ в (2.31) не пуст и либо в нем имеется вектор $d \neq 0$, либо $K^H(y^0) = \{0\}$.

Случай А. $K^H(y^0) = \{0\}$. В этой ситуации вектора $\nabla h_1(y^0), \nabla h_2(y^0), \dots, \nabla h_p(y^0)$ образуют базис в пространстве R^N (т.е. в этом случае, $p=N$), следовательно всегда

можно вектор $(\nabla f(y^0), v^0)$ разложить по нему с коэффициентами μ_1, \dots, μ_p , и положить $\lambda_i = 0$ ($i=1, \dots, m$).

Случай В. $\exists d \in K^H(y^0)$, что $d \neq 0$. При этом можно выделить два подслучая в (2.31): $K^Y \neq \emptyset$ и $K^Y = \emptyset$. Первый является основным, второй случай можно рассматривать во многом аналогично.

Лемма 2.2. При выполнении условий $K^f \cap K^Y = \emptyset$, $K^Y \neq \emptyset$ значение $d^0 = 0$ является слабо эффективным решением задачи

$$(\nabla f(y^0), d) \rightarrow \min, d \in \bar{K}^Y \quad (2.32)$$

Здесь черта означает замыкание множества.

Доказательство. Заметим, что $d^0 = 0$ — допустимая точка в (2.32).

Предположим, что она не слабо эффективна для этой задачи, тогда $\exists \eta \in \bar{K}^Y$, что $(\nabla f(y^0), \eta) < 0 \equiv (\nabla f(y^0), d^0)$. По предположению относительно $\eta : (\nabla g_j(y^0), \eta) \leq 0$ для $j \in J(y^0)$ и $(\nabla h(y^0), \eta) = 0$. Поскольку $K^Y \neq \emptyset$, то $\exists \xi \in K^Y : (\nabla g_j(y^0), \xi) < 0$ для $j \in J(y^0)$, $(\nabla h(y^0), \xi) = 0$. Образуя смесь $d^\varepsilon = \eta + \varepsilon \xi$ с малым $\varepsilon > 0$, получим вектор $d_\varepsilon \in K^f \cap K^Y$, что противоречит (2.31). Таким образом, лемма 2.2 доказана.

Воспользуемся доказанной леммой.

Задача (2.32) является выпуклой, а условие $K^Y \neq \emptyset$ гарантирует ее регулярность по достаточному условию Слейтера (теорема 2.3). Поэтому для нее можно записать необходимое и достаточное условие по теореме 2.5. Нетрудно видеть, что это условие совпадает с (2.24)–(2.26) и при этом окажется, что $v^0 \neq 0$, т.е. условие $K^Y \neq \emptyset$ является условием регулярности в исходной задаче.

Теперь рассмотрим оставшийся случай $K^Y = \emptyset$. Поскольку $K^H(y^0) \neq \emptyset$, то может быть лишь два случая. Либо $K_j^g(y^0) \cap K^H(y^0) = \emptyset \quad \forall j \in J(y^0)$, либо найдется такое $i < r = |J(y^0)|$ — мощности множества $J(y^0)$, что

$$K_{j_{i+1}}^g(y^0) \cap \dots \cap K_{j_r}^g(y^0) \cap K^H(y^0) \neq \emptyset, \text{ а } K_{j_i}^g(y^0) \cap \dots \cap K_{j_r}^g(y^0) \cap K^H(y^0) = \emptyset.$$

В первом из этих случаев вектор $\nabla g_{j_i}(y^0)$, очевидно, будет раскладываться по нормальям к равенствам, т.е. по системе $\nabla h_1(y^0), \dots, \nabla h_p(y^0)$, что позволит удовлетворить утверждению теоремы. Последний случай рассматривается аналогично случаю $K^Y \neq \emptyset$, но роль K^Y играет теперь непустое пересечение $\bigcap_{s=i+1}^r K_{j_s}^g(y^0) \cap K^H(y^0)$, а роль $K^f(y^0)$ — конус $K_{j_i}^g(y^0)$. Далее доказательство можно провести через лемму, аналогичную лемме 2.2. Теорема доказана.

Следствие 2.7.1 (достаточное условие регулярности для невыпуклых задач). Пусть в задаче FYE из (2.4), (2.5) $y^0 \in \text{int } E$, $f(y)$, $g(y)$ — дифференцируемы в точке y^0 , а $h(y)$ — непрерывно дифференцируема в y^0 . Если $\exists d \in R^N$, что $(\nabla h_s(y^0), d) = 0$ ($s = 1, \dots, p$) и $\forall j \in J(y^0) : (\nabla g_j(y^0), d) < 0$, при этом $\nabla h_1(y^0), \dots, \nabla h_p(y^0)$ -линейно независимы, то область Y регулярна в точке y^0 .

Этот факт непосредственно вытекает из хода доказательства теоремы.

Следствие 2.7.2 (вторая форма достаточного условия регулярности для невыпуклых задач). В задаче, приведенной в следствии 2.7.1, для регулярности области Y в точке $y^0 \in \text{int } E$ достаточно, чтобы набор векторов $\nabla h_1(y^0), \dots, \nabla h_p(y^0), \nabla g_{j_1}(y^0), \dots, \nabla g_{j_r}(y^0)$, где $\{j_1, \dots, j_r\} = J(y^0)$, был линейно независим.

Доказательство легко получается от противного.

2.2.2. Условия экстремума второго порядка

В теории математического программирования для гладких задач общего вида кроме дифференциальных условий первого порядка используются признаки второго порядка, позволяющие получить для локального экстремума определяющие его условия не только в необходимой, но и в достаточной форме. Эти результаты допускают частичное обобщение на многокритериальные задачи.

Нам предстоит разобраться с тем, как можно вывести такие обобщения и какие факты из однокритериальной оптимизации оказывается не переносимы на многокритериальный случай.

Следующий результат является прямым аналогом соответствующего факта для однокритериального случая.

Теорема 2.8. Пусть $f, g, h \in C^2(E)$ и для точки $y^0 \in Y$, являющейся внутренней для E ($y^0 \in \text{int } E$), выполнены условия A и B:

A. $\exists (v^0, \lambda^0, \mu^0), v^0 \geq 0, v^0 \neq 0, \lambda^0 \geq 0$, что $\lambda_j^0 g_j(y^0) = 0$ ($j=1, \dots, m$) и для функции Лагранжа (2.6) выполнено условие стационарности (2.17), т.е.

$$\nabla_y L(y^0, v^0, \lambda^0, \mu^0) = 0. \quad (2.33)$$

B. $\forall d \neq 0$ такого, что

$$(\nabla h(y^0), d) = 0, (\nabla g_j(y^0), d) = 0 \quad (j \in J(y^0)), \quad (2.34)$$

выполняется

$$d^T \Gamma_y^L(y^0, v^0, \lambda^0, \mu^0) d > 0. \quad (2.35)$$

Тогда при $n > 1$ точка y^0 – локально слабо эффективна, а при $n = 1$ является строгим локальным минимумом задачи.

В (2.35) используется Γ_y^L — матрица Гессе функции Лагранжа по переменным y .

ДОКАЗАТЕЛЬСТВО. Рассмотрим малую ε –окрестность точки y^0 . Возьмем произвольную точку $y \in O_\varepsilon(y^0) \cap Y$. Будем считать ε настолько малым, что множество номеров активных в точке y ограничений–неравенств $J(y) \subseteq J(y^0)$. Заметим, что $J(y)$ может отличаться от $J(y^0)$ за счет того, что точка y может быть выведена с границ части активных ограничений внутрь области Y .

Если вектора $d \neq 0$ и удовлетворяющие условиям пункта B существуют, то можно построить линейное, касательное в точке y^0 многообразию

$$M(y^0, y) = \{y^0 + d: (\nabla h(y^0), d) = 0, (\nabla g_j(y^0), d) = 0, j \in J(y)\}$$

к пересечению поверхностей, порожденных ограничениями–равенствами и активными в точке y неравенствами (рис.2.3).

Заметим, что возможен случай, когда $J(y) = \emptyset$ и ограничения–равенства отсутствуют. В этом случае следует считать, что $M(y^0, y) = R^N$, и все необходимые рассуждения следует провести с соответствующей поправкой.

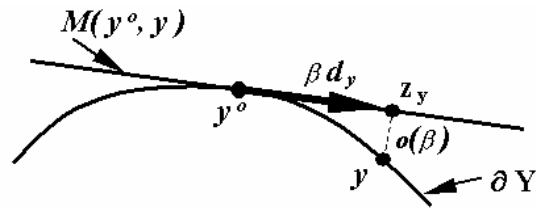


Рис.2.3. Фрагмент границы и касательное многообразие

Пусть z_y – проекция точки y на многообразие $M(y^0, y)$, запишем операцию проектирования как $z_y = \pi_{M(y^0, y)}(y)$. Введем d_y – нормированный вектор направления от y^0 к z_y :

$$d_y = (z_y - y^0) / \|z_y - y^0\|.$$

Тогда при некотором β : $z_y - y^0 = \beta d_y$. В силу достаточной гладкости функций h и g , $y - z_y = o(\beta)$. Это позволяет представить точку y в виде $y = y^0 + \beta d_y + o(\beta)$.

Построим оценки для скалярного произведения

$$\begin{aligned} (f(y), v^0) &\geq (f(y), v^0) + \underbrace{(g(y), \lambda^0)}_{\leq 0} + \underbrace{(h(y), \mu^0)}_{=0} = L(y, v^0, \lambda^0, \mu^0) = \\ &= L(y^0, v^0, \lambda^0, \mu^0) + \underbrace{(\nabla_y L(y^0, v^0, \lambda^0, \mu^0), y - y^0)}_{=0} + (d_y \beta + o(\beta))^T \Gamma_y^L(\xi, v^0, \lambda^0, \mu^0) (d_y \beta + o(\beta)) = \\ &= L(y^0, v^0, \lambda^0, \mu^0) + \beta^2 \underbrace{d_y^T \Gamma_y^L(\xi, v^0, \lambda^0, \mu^0) d_y}_{>0} + o(\beta^2) > L(y^0, v^0, \lambda^0, \mu^0) = (f(y^0), v^0) \quad (2.36) \end{aligned}$$

Итак, $\forall y \in Y \cap O_\varepsilon(y^0)$, $y \neq y^0$ при достаточно малом $\varepsilon > 0$ $(f(y), v^0) > (f(y^0), v^0)$. Отсюда при $n = 1$ (в силу $v^0 \neq 0$, $v^0 \geq 0$) следует, что y^0 – точка строгого локального минимума, а при $n > 1$ можно утверждать лишь (свойство 1.4 пункта 1.3.5), что y^0 – локально слабо эффективна. Теорема доказана.

Проанализируем теперь вопрос о необходимости условий (2.33)–(2.35) для слабой эффективности y^0 . Сформулируем следующую теорему.

Теорема 2.9. Пусть $f, g, h \in C^1(E)$, y^0 – локально слабо эффективная точка (локальный минимум при $n=1$), задача регулярна в точке y^0 и $y^0 \in \text{int } E$, тогда выполнено условие (A) теоремы 2.8. Если, кроме того, $f, g, h \in C^2(E)$ и задача – однокритериальна (т.е. $n=1$), то для всех векторов d , удовлетворяющих условиям (2.34) теоремы 2.8, выполняется неравенство

$$d^T \Gamma_y^L(y^0, v^0, \lambda^0, \mu^0) d \geq 0 \quad (2.37)$$

Замечание. Следует обратить внимание, что первое утверждение теоремы верно для любых задач, а второе – справедливо только для однокритериальных задач. Для многокритериального случая условие второго порядка не является необходимым.

Проследим ход рассуждений. Легко видеть, что первое утверждение следует из теоремы 2.7. Чтобы выяснить справедливость второго, предположим, что оно не выполняется и \exists вектор $d \neq 0$, удовлетворяющий (2.34), для которого неравенство (2.37) нарушается, т.е. $d^T \Gamma_y^L(y^0, v^0, \lambda^0, \mu^0) d < 0$. Введем коэффициент смещения β вдоль d . В силу достаточной гладкости g и h найдется поправка $o(\beta)$, что точка $y(\beta) = y^0 + \beta d + o(\beta)$ будет удовлетворять равенствам

$$h(y(\beta)) = 0 \text{ и } g_j(y(\beta)) = 0, j \in J(y^0).$$

Повторяя выкладки, аналогичные (2.36), получим

$$L(y(\beta), v^0, \lambda^0, \mu^0) < L(y^0, v^0, \lambda^0, \mu^0) = f(y^0, v^0)$$

Но

$$L(y(\beta), v^0, \lambda^0, \mu^0) = (f(y(\beta)), v^0) + \sum_{j \in J(y^0)} \lambda_j^0 \cdot g_j(y(\beta)) + (h(y(\beta)), \mu^0) = f(y(\beta), v^0)$$

в силу приведенных выше равенств, которым удовлетворяют в точке $y(\beta)$ функции g_j и h .

Следовательно, при малом β получаем, что $f(y(\beta), v^0) < (f(y^0), v^0)$ при $v^0 \geq 0, v^0 \neq 0$.

Рассмотрим это неравенство. При размерности векторного критерия $n=1$ мы получаем обычную однокритериальную задачу, в которой неравенство приведет к $f(y(\beta)) < f(y^0)$, что противоречит локальной оптимальности точки y^0 . Это означает, что свойство (2.37) выполняется для однокритериальных задач.

Пусть теперь $n > 1$. Приведем контрпример, показывающий, что рассматриваемое неравенство не противоречит локальной слабой эффективности точки y^0 и, следовательно, в многокритериальных задачах (2.37) может не выполняться. Структура примера поясняется на рис.2.4.

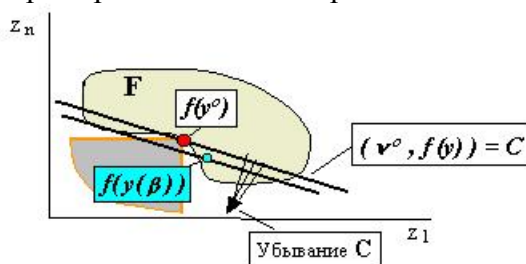


Рис.2. 4. Иллюстрация к контрпримеру

Видно, что при определенной структуре множества Слейтера и Парето существование в сколь угодно малой окрестности y^0 допустимой точки $y(\beta)$ со значением линейной свертки $\Psi_{v^0}(y(\beta)) < \Psi_{v^0}(y^0)$ не противоречит слабой эффективности точки y^0 .

Таким образом, неравенство (2.37) в многокритериальных задачах может нарушаться. Доказательство теоремы завершено.

2.3. Элементы теории двойственности в задачах математического программирования с одним критерием

Теория двойственности восходит к задачам линейного программирования, где она играет важную роль, являясь одним из инструментов трансформации постановки задачи к более удобной для решения форме.

В общих задачах нелинейного программирования роль теории двойственности более скромная, однако, она дает интересные подходы к некоторым методам учета ограничений в вычислительных процессах оптимизации.

Предыдущий материал по теории экстремума излагался для общей многокритериальной постановки и в значительной части представлял прямые обобщения результатов, известных для скалярных задач. В теории двойственности ситуация, к сожалению, иная. Ее обобщение на

многокритериальные задачи требует, как оказывается, привлечения новых конструкций по сравнению со скалярным случаем (см., например, главу 4 в [6]). Иными словами, возникают трудности при попытке такого изложения теории двойственности, когда результаты и форма их записи одинаково подходили бы для многокритериального и скалярного случаев.

По этой причине в настоящем курсе элементы теории двойственности будут представлены только для однокритериальных задач ((2.4),(2.5) при $n=1$).


Рассмотрим однокритериальную задачу, которую будем называть прямой задачей P .

Задача P

$$f(y) \rightarrow \min, y \in Y, f: Y \rightarrow R^1 \quad (2.38)$$

$$Y = \{y \in E: g(y) \leq 0, h(y) = 0\}, \quad (2.39)$$

где множество E , как это уже подчеркивалось в начале главы 2, может включать часть функциональных ограничений исходной постановки задачи.

 **Замечание.** Для однокритериальной задачи понятие слабо эффективного решения y^0 совпадает с понятием эффективного решения y^* , а также с понятием глобально-оптимального решения y^* . Поэтому далее для решения будет использоваться обозначение не y^0 , а y^* . Соответственно и значения множителей Лагранжа будут обозначаться как λ^*, μ^* .

Ограничимся в нашем рассмотрении задачами с выполненным условием регулярности (хотя бы для точки решения y^*). При этом множитель Лагранжа при $f(y)$ может быть принят равным $v^*=1$, поэтому, в отличие от общего случая (2.6), будем использовать функцию Лагранжа более простого вида

$$L(y, \lambda, \mu) = f(y) + (g(y), \lambda) + (h(y), \mu) \quad (2.40)$$

К постановке двойственной задачи проще всего придти, опираясь на полученные в разделе 2.1 условия экстремума. В теоремах 2.1, 2.4 было показано, что (применительно к рассматриваемой регулярной скалярной задаче) следующие условия являются достаточными для глобальной оптимальности точки $y^* \in Y$:

$$\exists \lambda^* \geq 0, \mu^*, \text{ что } \forall \lambda \geq 0 \in R^m, \mu \in R^p \text{ и } \forall y \in E \text{ выполняется} \\ L(y^*, \lambda, \mu) \leq L(y^*, \lambda^*, \mu^*) \leq L(y, \lambda^*, \mu^*) \quad (2.41)$$

Таким образом, в рассматриваемом классе задач седловая точка функции Лагранжа (2.40) по области $y \in E, \lambda \geq 0 \in R^m, \mu \in R^p$ определяет глобально-оптимальное решение y^* исходной задачи.

Структура выражения (2.41) подсказывает, что целесообразно ввести специальную функцию, определяющую нижнюю грань значений функции Лагранжа по $y \in E$. В отличие от (2.41) сделаем это при произвольных значениях множителей Лагранжа. А именно, для произвольных $\lambda \geq 0 \in R^m, \mu \in R^p$ введем функцию

$$L^*(\lambda, \mu) = \inf \{ L(y, \lambda, \mu) : y \in E \}, \quad (2.42)$$

которая при некоторых значениях параметров может принимать бесконечные значения равные $-\infty$.

Определение 2.3. Функцию $L^*(\lambda, \mu)$ называют функцией двойственной задачи.

Определение 2.4. Задачей двойственной по Лагранжу назовем задачу следующего вида (назовем ее задачей D)

Задача D

$$L^*(\lambda, \mu) \rightarrow \max, \lambda \geq 0 \in R^m, \mu \in R^p \quad (2.43)$$

Переменные λ и μ называют двойственными переменными.

С учетом (2.42) двойственная задача является *максиминной* задачей. Интуитивно кажется, что при достаточно общих условиях ее решение должно соответствовать точке λ^*, μ^*, y^* из (2.41) Несколько позднее этот вопрос будет рассмотрен точно.

Замечание. При записи задачи в постановке P ко множеству E может быть отнесена различная часть ограничений, что позволяет одну и ту же задачу записывать в разных формах P. Это приведет к неоднозначности формулировок двойственных задач D. Таким образом, одной задаче оптимизации можно сопоставить несколько форм записи P и несколько двойственных к ним задач D.

Теория двойственности устанавливает связи между решениями прямой и двойственной задач [1]. Первая (слабая) теорема двойственности устанавливает простое мажорантное соответствие между значениями $f(y)$ и $L^*(\lambda, \mu)$.

Теорема 2.10 (слабая теорема двойственности). $\forall y \in Y$ и $\forall (\lambda, \mu)$ с $\lambda \geq 0$ выполняется неравенство

$$f(y) \geq L^*(\lambda, \mu) \quad (2.44)$$

ДОКАЗАТЕЛЬСТВО вытекает из следующей простой оценки

$$L^*(\lambda, \mu) = \inf \{ f(y) + (\lambda, g(y)) + (\mu, h(y)) : y \in E \} \leq \underbrace{f(y)}_{\text{для } y \in Y} + \underbrace{(\lambda, g(y))}_{\leq 0} + \underbrace{(\mu, h(y))}_{= 0} \leq f(y).$$

В случае отсутствия ограничений–равенств и количестве ограничений–неравенств $m=1$ утверждение теоремы хорошо иллюстрирует рис.2.5.

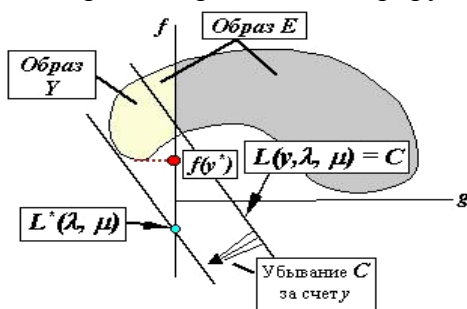


Рис.2. 5. Соотношение функций прямой и двойственной задач

Следствие 2.10.1.

A. $\inf \{ f(y) : y \in Y \} \geq \sup \{ L^*(\lambda, \mu) : \lambda \geq 0 \}$ (2.45)

B. Если нашлась точка $\bar{y} \in Y$ и $(\bar{\lambda}, \bar{\mu})$ с $\bar{\lambda} \geq 0$, что $f(\bar{y}) \leq L^*(\bar{\lambda}, \bar{\mu})$, то $\bar{y}, (\bar{\lambda}, \bar{\mu})$ являются решениями прямой и обратной задач P и D.

C. Если в прямой задаче P $\inf \{ f(y) : y \in Y \} = -\infty$, то $\forall (\lambda, \mu)$ с $\lambda \geq 0$ $L^*(\lambda, \mu) = -\infty$.

D. Если $\sup \{ L^*(\lambda, \mu) : \lambda \geq 0, \mu \in R^p \} = +\infty$, то $Y = \emptyset$.

В пояснении нуждается только последнее свойство. Его проще всего обосновать от противного. Допустим $\exists y \in Y$, тогда при $\lambda \geq 0$

$$L^*(\lambda, \mu) \leq f(y) + \underbrace{(g(y), \lambda)}_{\leq 0} + \underbrace{(h(y), \mu)}_{= 0} \leq f(y) < +\infty,$$

что противоречит заданному условию.

Центральный факт из теории двойственности – сильная теорема двойственности. Она показывает, что при условиях выпуклости и регулярности оптимальные значения функций в задачах P и D совпадают. Этот факт вполне понятен в рамках геометрической интерпретации задач на последнем рисунке.

Теорема 2.11 (сильная теорема двойственности). Пусть множество E выпукло и не пусто, f и g – выпуклы на E , $h(y)$ – аффинная ($h(y) = Ay - b$). Пусть, кроме того, выполнены условия регулярности Слейтера (теорема 2.3): $\exists \bar{y} \in E : h(\bar{y}) = 0, g(\bar{y}) < 0$ и, кроме того, $0 \in \text{int } h(E)$, тогда

$$\inf \{ f(y) : y \in Y \} = \sup \{ L^*(\lambda, \mu) : \lambda \geq 0 \} \quad (2.46)$$

Если \inf конечен, то \sup достигается в (λ^*, μ^*) с $\lambda^* \geq 0$.

Если \inf достигается в точке y^* , то $(\lambda^*, g(y^*)) = 0$.

ДОКАЗАТЕЛЬСТВО.

1. Пусть $\inf = -\infty$, тогда по следствию 2.10.1(C) из слабой теоремы двойственности $\sup = -\infty$, т.е. равенство (2.46) будет выполнено.

2. Пусть \inf достигается в некоторой точке $y^* \in Y$, тогда по теоремам 2.2–2.3 $\exists (\lambda^*, \mu^*)$ с $\lambda^* \geq 0$, что $(\lambda^*, g(y^*)) = 0$ и

$$L(y^*, \lambda^*, \mu^*) = \min \{ L(y, \lambda^*, \mu^*) : y \in Y \} = L^*(\lambda^*, \mu^*).$$

Но значение в левой части равенства равно $f(y^*) + (\lambda^*, g(y^*)) + (\mu^*, h(y^*))$, где два последних члена равны нулю, следовательно, $f(y^*) \leq L^*(\lambda^*, \mu^*)$.

Тогда по следствию 2.10.1(B) из слабой теоремы двойственности получаем выполнение (2.46), также окажется выполненным и условие $(\lambda^*, g(y^*)) = 0$.

Остался не рассмотренным случай, когда $\inf \{ f(y) : y \in Y \} = C \neq -\infty$, но он не достигается. Тогда $\forall y \in E$ система

$$\begin{cases} g(y) \leq 0, h(y) = 0 \\ f(y) \leq C \end{cases}$$

не совместна.

Если рассмотреть множество $S = \bigcup_{y \in E} S(y)$ из доказательства теоремы 2.2, то точка $(C, 0, 0)$ будет его граничной точкой. Тогда по теореме об отделимости выпуклого множества S и его граничной точки $(C, 0, 0)$ $\exists (v^*, \lambda^*, \mu^*) \neq 0$ что $\forall \varepsilon \geq 0, \delta \geq 0, \forall y \in E$:

$$v^*(f(y) + \varepsilon) + (\lambda^*, g(y) + \delta) + (\mu^*, h(y)) \geq C \cdot v^*$$

Далее, аналогично доказательству теоремы 2.2, можно показать, что $v^* \geq 0, \lambda^* \geq 0$, где, в силу регулярности задачи, $v^* \neq 0$ и его можно положить равным 1. Выбирая в неравенстве $\varepsilon = 0, \delta = 0$, получим, что $L(y, v^*, \lambda^*, \mu^*) \geq C$, поэтому $L^*(\lambda^*, \mu^*) \geq C$. По следствию 2.10.1(A) это доказывает справедливость утверждения теоремы.

Рассмотрим применение теории двойственности в аспекте построения вычислительных методов решения задач с ограничениями. Понятно, что в

условиях теоремы 2.11 при существовании конечного глобального минимума задача P можно определить его, решая двойственную задачу D в пространстве $(\lambda, \mu) \in R^{m+p}$ с простыми ограничениями $\lambda \geq 0$. Проблема заключается в том, что функция двойственной задачи $L^*(\lambda, \mu)$ определяется через решение вспомогательной экстремальной задачи (2.42), где минимум ищется только по области E , в качестве которой обычно используется либо все пространство R^N , либо область (2.3) в виде N -мерного параллелепипеда.

Таким образом, для выпуклой регулярной задачи

$$f(y^*) = \max_{\mu \in R^p; \lambda \in R^m; \lambda \geq 0} \overbrace{\min_{y \in E} L(y, \lambda, \mu)}^{L^*(\lambda, \mu)} \quad (2.47)$$

Выигрыш этого подхода состоит в том, что нигде не требуется решать экстремальные задачи со сложными функциональными ограничениями, а трудности вычислительной реализации связаны с неявным заданием максимизируемой функции $L^*(\lambda, \mu)$ внешней экстремальной задачи.

Дальнейший материал этого раздела направлен на изучение свойств функции $L^*(\lambda, \mu)$. Сейчас будет показано, что она обладает достаточно хорошей структурой.

Теорема 2.12. Пусть E -компакт, а $f, g, h \in C(E)$, т.е. непрерывны на E . Тогда функция $L^*(\lambda, \mu)$ двойственной задачи вогнута на множестве $\lambda \geq 0 \in R^m, \mu \in R^p$.

ДОКАЗАТЕЛЬСТВО. Выберем два допустимых значения для λ и μ : (λ^1, μ^1) , (λ^2, μ^2) и произвольное $\alpha \in (0, 1)$. Сделаем оценку

$$\begin{aligned} L^*(\lambda^\alpha, \mu^\alpha) &= \inf \{ f(y) + (\lambda^\alpha, g(y)) + (\mu^\alpha, h(y)) : y \in E \} = \\ &= \inf \{ \alpha L(y, \lambda^1, \mu^1) + (1-\alpha)L(y, \lambda^2, \mu^2) : y \in E \} \geq \alpha L^*(\lambda^1, \mu^1) + (1-\alpha)L^*(\lambda^2, \mu^2). \end{aligned}$$

Вогнутость доказана.

При построении большинства вычислительных методов оптимизации требуется измерение градиента целевой функции. Возникает вопрос о дифференцируемости $L^*(\lambda, \mu)$ по (λ, μ) , а также о способе вычисления ее градиента.

Приведем без доказательства² теорему, дающую ответ на этот вопрос.

Теорема 2.13. Если E – компакт, функции $f, g, h \in C(E)$, и $\forall (\lambda, \mu)$ с $\lambda \geq 0$ точка глобального минимума $y^*(\lambda, \mu)$ функции Лагранжа по y на E единственна, то функция $L^*(\lambda, \mu)$ дифференцируема в указанной области по λ, μ и

$$\nabla_{\lambda, \mu} L^*(\lambda, \mu) = (g(y^*(\lambda, \mu)); h(y^*(\lambda, \mu)))^T \quad (2.48)$$

При нарушении единственности точки глобального минимума $y^*(\lambda, \mu)$ функция $L^*(\lambda, \mu)$ может оказаться не дифференцируемой, однако в этом случае можно воспользоваться понятием субградиента вогнутой функции представляющего направления ее локального роста.

² Обоснование аналогичного свойства можно найти, например, в книге М.Базара, К.Шетти [1], раздел 6.3

Определение 2.5. Субградиентом вогнутой функции $f(y)$ в точке \bar{y} называют любое направление d , при котором $\forall y \in E$ выполняется неравенство $f(y) \leq f(\bar{y}) + (d, y - \bar{y})$.

Это определение полезно сопоставить со свойством 1.7 для выпуклых функций из раздела 1.4.

Геометрически субградиенты вогнутой функции f представляют соответствующие нормали к гиперплоскостям, являющимся опорными в точке \bar{y} ко множеству

$$Y(\bar{y}) = \{y \in R^N : f(y) > f(\bar{y})\},$$

что иллюстрирует рис.2.6. На нем представлены векторы $-sub\nabla f(\bar{y})$, обратные по направлению к субградиентам.

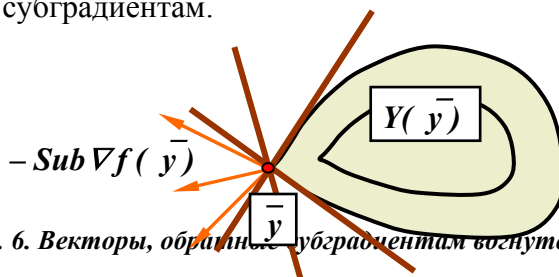


Рис.2. 6. Векторы, обратные субградиентам вогнутой функции

Имеет место следующая теорема, обобщающая предыдущую (ее обоснование приведено в той же книге М.Базара, К.Шетти [1]).

Теорема 2.14. Если E -компакт, функции $f, g, h \in C(E)$, то вектор

$$sub\nabla L^*(\lambda, \mu) = (g(y^*(\lambda, \mu)); h(y^*(\lambda, \mu)))^T \in Sub\nabla L^*(\lambda, \mu), \quad (2.49)$$

т.е. является одним из субградиентов функции $L^*(\lambda, \mu)$.

Существование простых формул для вычисления градиента или субградиента функции двойственной задачи представляется достаточно удивительным. Покажем как формула (2.48) получается в наиболее простом случае, когда нижняя грань в (2.42) достигается во внутренней точке множества E , т.е. $y^*(\lambda, \mu) \in int E$.

По определению имеем

$$L^*(\lambda, \mu) = f(y^*(\lambda, \mu)) + \sum_{i=1}^m \lambda_i g_i(y^*(\lambda, \mu)) + \sum_{j=1}^p \mu_j h_j(y^*(\lambda, \mu)).$$

Для внутренней точки минимума

$$\nabla_{y^*} L(y^*, \lambda, \mu) = \nabla f(y^*(\lambda, \mu)) + \sum_{i=1}^m \lambda_i \nabla g_i(y^*(\lambda, \mu)) + \sum_{j=1}^p \mu_j \nabla h_j(y^*(\lambda, \mu)) \equiv 0.$$

Это приводит к обнулению последнего произведения в нижеследующем выражении для частных производных двойственной функции Лагранжа:

$$\begin{aligned} \frac{\partial L^*(\lambda, \mu)}{\partial \lambda_s} &= \left(\frac{\partial y^*(\lambda, \mu)}{\partial \lambda_s} \right)^T \cdot \nabla f(y^*(\lambda, \mu)) + g_s(y^*(\lambda, \mu)) + \\ &+ \sum_{i=1}^m \lambda_i \left(\frac{\partial y^*(\lambda, \mu)}{\partial \lambda_s} \right)^T \nabla g_i(y^*(\lambda, \mu)) + \sum_{j=1}^p \mu_j \left(\frac{\partial y^*(\lambda, \mu)}{\partial \lambda_s} \right)^T \nabla h_j(y^*(\lambda, \mu)) = \end{aligned}$$

$$= g_s(y^*(\lambda, \mu)) + \left(\frac{\partial y^*(\lambda, \mu)}{\partial \lambda_s} \right)^T \cdot \nabla_y L(y^*(\lambda, \mu), \lambda, \mu) = g_s(y^*(\lambda, \mu)).$$

Аналогично можно показать, что

$$\frac{\partial L^*(\lambda, \mu)}{\partial \mu_i} = h_i(y^*(\lambda, \mu)).$$

Выполненные построения приводят к следующему алгоритму решения регулярных выпуклых задач с ограничениями (см. условия теоремы 2.11) через решение задачи (2.47):

$$y^{k+1} = \arg \min \{ L(y, \lambda^k, \mu^k) : y \in E \}. \quad (2.50)$$

$$\lambda^{k+1} = (\lambda^k + x^k \bar{g}(y^{k+1}))_+ \quad (2.51)$$

$$\mu^{k+1} = \mu^k + x^k h(y^{k+1}), \quad (2.52)$$

где (2.50) приводит к вычислению значения функции двойственной задачи $L^*(\lambda^k, \mu^k)$, а в (2.51), (2.52) выполняются шаги ее максимизации в направлении субградиента. Операция проектирования $(\cdot)_+$ обеспечивает выполнение условия $\lambda^{k+1} \geq 0$, а замена $g(y)$ на $\bar{g}(y)$ выполняется по правилу

$$\bar{g}_i(y^{k+1}) = \begin{cases} 0, & \text{если } \lambda_i^k = 0, g_i(y^{k+1}) \leq 0 \\ g_i(y^{k+1}), & \text{иначе} \end{cases}.$$

Коэффициент $x^k \in R^l$, определяющий величину смещения вдоль выбранного направления, может выбираться с помощью специальной вычислительной процедуры (которая будет рассмотрена значительно позднее — в главе 7 при изучении методов локальной оптимизации) из условия приближенного достижения максимума в задаче

$$\max \{ L(y^{k+1}, \lambda^k + x \bar{g}(y^{k+1}), \mu^k + x h(y^{k+1})) : \lambda^k + x \bar{g}(y^{k+1}) \geq 0, x \geq 0 \}.$$

Вычислительные методы рассмотренного типа будут обсуждаться в главе 3.

Лист регистрации изменений

Дата	Автор	Комментарии
06.06.02	Городецкий С.Ю.	Создание документа
14.07.02	Городецкий С.Ю.	Внесение изменений
15.07.02	Городецкий С.Ю.	Внесение изменений
16.07.02	Городецкий С.Ю.	Внесение изменений
25.12.02	Городецкий С.Ю.	Окончательная редакция версии 1
10.09.03	Городецкий С.Ю.	Изменения для версии 2
13.09.03	Городецкий С.Ю.	Окончательная редакция версии 2

Глава 3. Общие методы учета ограничений в задачах математического программирования

Ограничения на переменные в задачах оптимизации можно разделить на две группы: *специальные ограничения* и *ограничения общего вида*. Специальными называют такие ограничения, для учета которых существуют и применяются в программной системе особые приемы и алгоритмы. К специальным ограничениям в первую очередь следует отнести линейные ограничения. В простейшем и, одновременно, наиболее распространенном случае линейные ограничения–неравенства присутствуют в виде двусторонних ограничений на переменные. Специальные методы учета линейных ограничений будут рассмотрены позднее — в главе 7, посвященной вопросам локальной оптимизации.

Все те ограничения, для учета которых специальные алгоритмы не применяются, называют ограничениями общего вида. Если, например, специфика линейного ограничения не будет специально учитываться методом оптимизации, то это ограничение следует рассматривать как ограничение общего вида.

Ограничения общего вида наиболее просто учитываются с помощью сведения задачи с такими ограничениями к одной или последовательности задач, в которых подобные ограничения отсутствуют. Это достигается за счет использования вспомогательных задач, минимизируемые функции которых строятся с учетом не только целевой функции исходной задачи, но и функций всех общих ограничений. Методы построения таких вспомогательных задач рассматриваются в этой главе.

3.1. Общие методы учета ограничений, обзор методов

Рассмотрим задачу математического программирования

$$f(y) \rightarrow \min, y \in Y, f: Y \rightarrow R^1, \quad (3.1)$$

$$Y = \{y \in E : g(y) \leq 0, h(y) = 0\}, \quad (3.2)$$

в предположении, что функции g и h определяют ограничения общего вида, а все специально учитываемые ограничения отнесены к области E . Этот момент отличает постановку (3.1), (3.2) от внешне похожей задачи (2.38), (2.39).

Предположим, что в распоряжении вычислителя имеется метод поиска минимума функций в области E без дополнительных функциональных ограничений. Рассмотрим следующий вопрос: как с использованием такого вычислительного метода найти (хотя бы приближенно) решение задачи с дополнительными функциональными ограничениями (3.2)?

Как уже отмечалось, основная идея общих методов учета ограничений состоит в том, чтобы для задачи (3.1), (3.2) с функциональными ограничениями построить специальное семейство задач без функциональных ограничений

$$S_\beta(y) \rightarrow \min, y \in E, \quad (3.3)$$

в которых параметры β должны подстраиваться таким образом, чтобы предельные точки последовательности решений y_β^* этих задач являлись решениями исходной задачи. Целевые функции $S_\beta(y)$ вспомогательных задач (3.3) строятся по функциям f, g и h исходной задачи.

Методы такого типа называют общими потому, что они не требуют разработки новых вычислительных процедур условной оптимизации, а опираются на использование лишь общих методов поиска минимума в областях простой структуры при специальных способах построения последовательности вспомогательных задач (3.3). Заметим, что в зависимости от вида этих задач, коэффициент β может подстраиваться различным образом: либо изменяться после решения каждой вспомогательной задачи, либо — непосредственно в процессе их решения.

Известно несколько подходов к реализации этой общей идеи. Они представлены на схеме рис. 3.1.

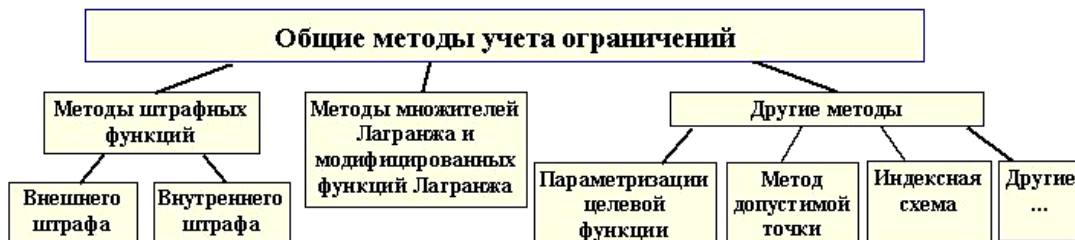


Рис.3. 1. Классификация общих методов учета ограничений

Основные из этих методов будут рассмотрены далее в этой главе.

3.2. Метод внешнего штрафа

3.2.1. Общее описание и некоторые свойства

Семейство методов штрафных функций представлено в этом разделе методом внешнего штрафа. Он широко используется при практических расчетах, особенно при решении задач локальной оптимизации. Его следует признать наиболее наглядным из общих методов учета ограничений. Информацию по методу внутреннего штрафа (метод барьеров) можно найти, например, в книге М.Базара, К.Шетти [1], раздел 9.3.

Введем понятие функции штрафа.

Определение 3.1. Непрерывную на E функцию $H(y)$, удовлетворяющую условию

$$H(y) = \begin{cases} 0, & y \in Y \\ > 0, & y \in E \setminus Y, \end{cases} \quad (3.4)$$

называют функцией (внешнего) штрафа.

Поскольку геометрическая структура области Y не известна и косвенно определяется через функции ограничений в (3.2), то штраф $H(y)$ можно определить в виде функции, зависящей от значений $g(y)$ и $h(y)$. Существуют различные способы задания этой функции (см., например, раздел 6.5 в книге [8]). Наиболее часто используется *степенная функция штрафа* вида

$$H(y) = \sum_{j=1}^m \left(\max \{0; c_j g_j(y)\} \right)^r + \sum_{s=1}^p |C_s \cdot h_s(y)|^r \quad (3.5)$$

где $c_1, \dots, c_m > 0, r > 0$.

Если считать, что функции g и h уже приведены к единому масштабу общей физической размерности за счет ранее выполненной нормировки, то

коэффициенты c_i и C_j можно опустить. Функция $H(y)$ трактуется как “штраф”, который накладывается на целевую функцию за нарушение ограничений.

Целевая функция задачи со штрафом строится в виде

$$S_\gamma(y) = f(y) + \gamma H(y) \quad (3.6)$$

с $\gamma > 0$. Параметр γ называют *коэффициентом штрафа*. Метод заключается в решении последовательности задач со штрафом вида (3.3) при $\beta = \gamma_k$ для возрастающей последовательности $\gamma = \gamma_k$


$$S_\gamma(y) \rightarrow \min, y \in E \quad (3.7)$$


Укажем на некоторые очевидные свойства метода штрафа.

Свойства 3.1.

- A. При $r \leq 1$ функция $H(y)$ из (3.5) не является гладкой по переменным g и h .*
- B. При $r > 1$ $H(y)$ становится непрерывно дифференцируемой, а при $r > 2$ – дважды непрерывно дифференцируемой по g и h .*
- C. При больших значениях γ функции задачи со штрафом $S_\gamma(y)$ становятся сильно овражными.*

Под *овражностью* здесь понимается существование многообразий в пространстве переменных y , вдоль которых функция $S_\gamma(y)$ изменяется много медленнее, чем при смещении вдоль направлений, локально ортогональных к ним.

 **Замечание 1.** Относительно свойств *A* и *B* следует иметь в виду, что большинство методов локальной оптимизации, часто используемых в сочетании с методом штрафов, весьма чувствительны к степени гладкости и нарушению гладкости.

 **Замечание 2.** Если функция штрафа (3.5) является гладкой, то нельзя гарантировать, что решение задачи со штрафом (3.6), (3.7) при каком-либо конечном γ будет совпадать с решением исходной задачи (3.1)–(3.2). Поэтому для гладкой функции штрафа приходится рассматривать бесконечно возрастающую последовательность значений коэффициента штрафа.

Для обоснования этого утверждения достаточно привести пример задач, в которых наблюдается описанная ситуация. Пусть функция штрафа достаточно гладкая, а целевая функция исходной задачи строго убывает в некоторой граничной точке y допустимой области Y в направлении d , выводящем из множества Y и не выводящем из E . Предположим также, что скорость убывания отделена от нуля. Пусть, кроме того, в этой точке достигается минимум исходной задачи. Поскольку в допустимых точках функция штрафа тождественно равна нулю, то в граничной для Y точке y градиент гладкой функции штрафа обратится в ноль. Следовательно производная штрафной добавки, вычисленная в точке y в направлении d будет равна нулю вне зависимости от значения коэффициента штрафа. Из этого следует, что во вспомогательных задачах со штрафом функции $S_\gamma(y)$ будут локально строго убывать в точке y в направлении d , поэтому точка y не будет являться решением вспомогательных задач со штрафом ни при каких конечных значениях коэффициента γ .

На рис. 3.2 показано поведение функций задачи со штрафом в случае одного переменного и одного ограничения $g(y) \leq 0$ при показателе степени в штрафе $r=2$. Можно видеть изменение функции штрафа при увеличении коэффициента γ , а также вид функции штрафной задачи $S_\gamma(y)$ при одном из значений коэффициента штрафа.

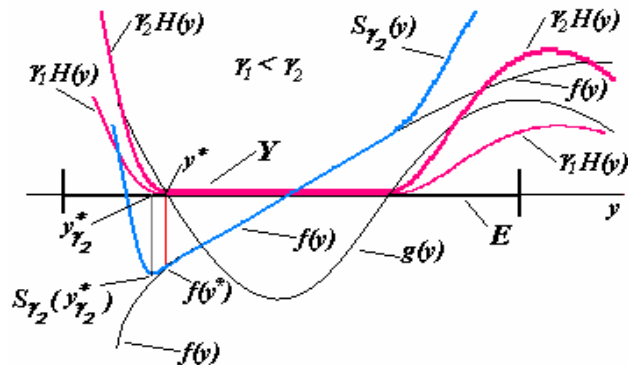


Рис.3. 2. Пример поведения функции штрафной задачи

Следует обратить внимание на то, что в задаче, представленной на рис.3.2, за счет дифференцируемости функции штрафа $H(y)$, ни при каком конечном значении γ точка минимума функции $S_\gamma(y)$ не будет совпадать с решением исходной задачи.

При достаточно общих условиях можно обосновать сходимость процедуры метода штрафов при $\gamma \rightarrow \infty$.

3.2.2. Исследование сходимости и алгоритм настройки коэффициента штрафа

Изучая вопросы сходимости метода штрафа следует учитывать, что решение задач (3.7) при $\gamma = \gamma_k$ всегда происходит с некоторой вычислительной погрешностью $\varepsilon = \varepsilon_k$. Приближенные решения y_γ^* удовлетворяют условию

$$S_\gamma(y_\gamma^*(\varepsilon)) \leq S_\gamma(y_\gamma^*) + \varepsilon \quad (3.8)$$

Теорема 3.1 (о достаточных условиях сходимости метода штрафов). Пусть выполнены следующие условия:

1. Функции $f(y)$, $g(y)$, $h(y)$ и $H(y)$ непрерывны на E ;
2. Существуют решения $y^* \in Y^*$ исходной задачи, а при достаточно больших $\gamma > 0$ — решения y_γ^* вспомогательных задач со штрафом;
3. Существует компакт $K \subseteq E$, что для достаточно больших γ и малых ε точки $y_\gamma^*(\varepsilon)$, определяющие приближенные решения задач со штрафом (согласно (3.8)), содержатся в K ;
4. Используются значения $\gamma = \gamma_k \rightarrow \infty$, $\varepsilon = \varepsilon_k \rightarrow 0$ при $k \rightarrow \infty$, $\gamma_k > 0, \varepsilon_k \geq 0$.

Тогда предельные точки y_∞^* последовательности $y_\gamma^*(\varepsilon)$ являются решениями исходной задачи и существуют пределы

$$\lim_{k \rightarrow \infty} f(y_{\gamma_k}^*(\varepsilon_k)) = f(y^*) \quad (3.9)$$

$$\lim_{k \rightarrow \infty} \rho(y_{\gamma_k}^*(\varepsilon_k), Y^*) = 0 \quad (3.10)$$

Доказательство. Очевидно, что при существовании решений исходной и вспомогательных задач выполняется неравенство

$$\begin{aligned} f(y^*) + \varepsilon_k &= S_{\gamma_k}(y^*) + \varepsilon_k \geq S_{\gamma_k}(y_{\gamma_k}^*) + \varepsilon_k \geq S_{\gamma_k}(y_{\gamma_k}^*(\varepsilon_k)) = \\ &= f(y_{\gamma_k}^*(\varepsilon_k)) + \gamma_k H(y_{\gamma_k}^*(\varepsilon_k)). \end{aligned}$$

Поделим его на γ_k , тогда

$$f(y^*)/\gamma_k + \varepsilon_k/\gamma_k \geq f(y_{\gamma_k}^*(\varepsilon_k))/\gamma_k + H(y_{\gamma_k}^*(\varepsilon_k)).$$

Рассмотрим произвольную предельную точку y_∞^* последовательности $y_{\gamma_k}^*(\varepsilon_k)$ и перейдем к пределу на сходящейся к ней подпоследовательности (эти действия обоснованы в силу принадлежности точек $y_{\gamma_k}^*(\varepsilon_k)$ компакту K). В пределе из предыдущего неравенства получим $0 \geq H(y_\infty^*) \geq 0$, следовательно, y_∞^* – допустима.

Кроме того, используя полученную выше оценку, имеем

$$f(y_{\gamma_k}^*(\varepsilon_k)) \leq f(y_{\gamma_k}^*(\varepsilon_k)) + \gamma_k H(y_{\gamma_k}^*(\varepsilon_k)) \leq f(y^*) + \varepsilon_k.$$

Переходя к пределу на сходящейся к y_∞^* подпоследовательности получим, что $f(y_\infty^*) \leq f(y^*)$. В силу допустимости предельной точки $f(y_\infty^*) = f(y^*)$. Следовательно, $y_\infty^* \in Y^*$.

Докажем теперь существование пределов (3.9), (3.10) для всей последовательности $y_{\gamma_k}^*(\varepsilon_k)$. Предположим, что предел (3.9) не существует, тогда найдется такая подпоследовательность $k=k_s$, что $\forall k = k_s, s=1,2,\dots$

$$|f(y_{\gamma_k}^*(\varepsilon_k)) - f(y^*)| > \delta > 0 \quad (3.11)$$

Однако, в силу доказанного ранее, из $y_{\gamma_k}^*(\varepsilon_k)$, ($k=k_s$) можно выделить подпоследовательность, сходящуюся к одному из решений y^* , что противоречит (3.11). Таким образом, соотношение (3.9) доказано.

Используя тот же прием в рассуждениях, можно обосновать (3.10). Теорема доказана.

Для вычислительной реализации метода штрафов необходимо выбрать алгоритм, определяющий закон изменения коэффициента штрафа и точности решения штрафных задач. Опишем один из возможных алгоритмов. Он основан на том, что контролируется убывание невязки по ограничениям на каждом шаге.

Невязкой в точке y назовем величину $G(y)$, показывающую степень нарушения ограничений в этой точке. Определим невязку с учетом нормировочных коэффициентов $c_j > 0$, использованных в функции штрафа

$$G(y) = \max \{ \max \{ c_j g_j(y) : j=1, \dots, m \}; 0; \max \{ |C_s h_s(y)| : s=1, \dots, p \} \}, \quad (3.12)$$

ОПИСАНИЕ АЛГОРИТМА настройки параметров метода штрафов.

ШАГ 0. Задаются: $\varepsilon_0 > 0$ – начальная точность решения штрафных задач, $\varepsilon > 0$ – требуемая точность их решения, $\delta > 0$ – требуемая точность по ограничениям, $0 < \alpha < 1$ – ожидаемый коэффициент убывания невязки по ограничениям на шаге, $\beta > 1$ – коэффициент увеличения штрафа, $\beta_1 > 1$ – дополнительный коэффициент увеличения штрафа, $0 < \nu < 1$ – коэффициент повышения точности решения штрафной задачи, $\gamma = \gamma_0$ – начальное значение коэффициента штрафа.

ШАГ 1. Решаем задачу со штрафом (3.7), (3.6) с точностью ε_0 , получаем оценку решения $y_{\gamma_0}^*(\varepsilon_0)$. Вычисляем начальную невязку полученного решения по

ограничениям $G_0 = G(y_{\gamma_0}^*(\varepsilon_0))$. Полагаем $k=0$ – номер выполненной итерации. Строим увеличенное значение коэффициента штрафа $\gamma_{k+1} = \beta \gamma_k$ и изменяем значение точности $\varepsilon_{k+1} = \nu \varepsilon_k$.

ШАГ 2. Проверяем критерий останова: если $G_k < \delta$ и точность решения задачи $\varepsilon_k \leq \varepsilon$, то выполняем останов процесса решения. Если $G_k < \delta$, но точность решения задачи $\varepsilon_k > \varepsilon$, то полагаем $\varepsilon_{k+1} = \varepsilon$ и переходим на шаг 3, если же $G_k \geq \delta$ — сразу переходим на шаг 3.

ШАГ 3. Решаем задачу со штрафом (3.7), (3.6) при $\gamma_k = \gamma_{k+1}$ с точностью ε_{k+1} , получаем оценку решения $y_{\gamma_{k+1}}^*(\varepsilon_{k+1})$. Вычисляем невязку полученного решения по ограничениям $G_{k+1} = G(y_{\gamma_{k+1}}^*(\varepsilon_{k+1}))$.

ШАГ 4. Если $G_{k+1} < \alpha G_k$, то полагаем $\gamma_{k+2} = \beta \gamma_{k+1}$, иначе $\gamma_{k+2} = \beta_1 \beta \gamma_{k+1}$. Повышаем точность $\varepsilon_{k+2} = \nu \varepsilon_{k+1}$, Полагаем $k=k+1$. Возвращаемся на шаг 2.

Описанный алгоритм формально обеспечивает выполнение требований теоремы 3.1, однако необходимо иметь в виду, что решение задачи с заданной точностью не гарантируется методами локального уточнения решений, которые обычно используются на шаге 3. Кроме того, при повышении требований к точности в малой окрестности решения начинают сказываться ошибки конечноразрядной арифметики. Поэтому при практических расчетах не следует уменьшать ε_k более некоторого порогового значения.

При проведении практических расчетов необходимо также учитывать возможные грубые ошибки в работе вычислительных методов, в результате которых вместо определения глобального минимума задач со штрафом происходит определение их локального минимума не являющегося глобальным. В этом случае может оказаться, что при увеличении коэффициента штрафа γ_k невязка по ограничениям в новой штрафной задаче не будет уменьшена. В этом случае необходимо прервать вычисления и попытаться подобрать другой вычислительный метод для поиска решения штрафных задач.

3.2.3. Структура возникающих задач со штрафом и характер приближения оценок к решению

Для того, чтобы можно было прогнозировать характер поведения вычислительных методов при решении вспомогательных задач со штрафом (3.7), необходимо изучить характерные особенности, имеющиеся в структуре $S_{\gamma_k}(y)$ — функций штрафных задач. В начале данного пункта проведем неформальное обсуждение возникающих вычислительных особенностей. Эти особенности могут быть обусловлены тремя причинами.

Первая связана с использованием больших значений коэффициента штрафа. Это приводит к тому, что в задачах, имеющих решение на границе допустимой области Y , будут возникать функции $S_{\gamma_k}(y)$ с сильно овражной структурой.

Причина заключается в том, что минимизируемая функция $f(y)$ исходной задачи, в общем случае, изменяется вдоль границы области относительно медленно, а функция штрафа постоянна: $H(y)=0$. Если же точка начинает удаляться от границы допустимой области в область недопустимости, то, за счет больших значений коэффициента штрафа, функция $S_{\gamma_k}(y)$ будет быстро изменяться. Таким образом, у этой функции наблюдается «овраг».

На рис.3.3 приведен пример, показывающий изменение изолиний функции $S_{\gamma}(y)$ штрафной задачи вида (3.5)–(3.7), возникающей при поиске минимума функции $(y_1-0,1)^2+0,1(y_2-0,8)^2$ с ограничением $-9y_1^2-y_2^2 \leq -1$ при увеличении коэффициента штрафа со значения $\gamma=1$ до значения $\gamma=20$. В функции штрафа использован показатель степени $r=2$. Допустимая область выделена более темным цветом.

Вторая особенность связана с тем, что функция штрафа $H(y)$ из (3.5), аддитивно входящая в функцию $S_{\gamma_k}(y)$, может иметь нарушение гладкости, начиная с некоторого порядка, даже при гладких функциях ограничений–неравенств. Например, у функции штрафа вида (3.5) при показателе степени $0 < r \leq 1$ на границе нарушения ограничений все частные производные, в общем случае, будут разрывны, а при $1 < r \leq 2$ первые производные станут непрерывными, но производные более высокого порядка будут терпеть разрыв. Нарушение гладкости может отрицательно сказываться на работе методов локальной оптимизации при решении задач методом штрафных функций. В то же время, повышение гладкости штрафа ухудшает скорость сходимости метода штрафов за счет того, что вблизи границы области Y штрафная добавка становится бесконечно малой более высокого порядка, чем расстояние до границы допустимой области. Ситуация усложняется тем, что именно в окрестности этой границы (если решение исходной задачи (3.1)–(3.2) лежит на границе области Y) функция $S_{\gamma_k}(y)$ сильно «овражна» при больших значениях γ_k . Таким образом, при выборе вычислительного метода оптимизации, используемого в методе штрафов, необходимо учитывать возможное нарушение гладкости вблизи дна оврага, вдоль которого обычно выполняется поиск минимума. Заметим, что приведенные здесь рассуждения, носящие качественный характер, можно подкрепить точными оценками скорости сходимости метода штрафов, увязав их со значением показателя степени r в функции штрафа. Эти оценки приводятся далее в теореме 3.2.

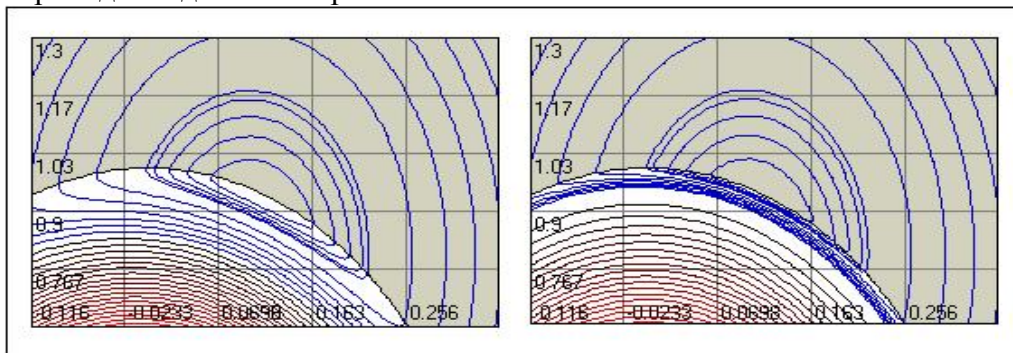



Рис.3.3. Увеличение овражности функции задачи со штрафом $S_{\gamma}(y)$ при возрастании коэффициента штрафа γ с 1 до 20

Третья особенность функции $S_{\gamma_k}(y)$ связана с тем, что порядок ее роста при отклонениях от границы области Y может быть существенно различен, в зависимости от того, происходит отклонение внутрь области Y или вне ее. Эта особенность может оказаться существенной для большой группы методов, основанных на квадратичной модели минимизируемой функции (модифицированный метод Ньютона, квазиньютоновские методы, метод сопряженных градиентов).

 **Замечание.** Таким образом, функции штрафных задач обладают характерными особенностями, ухудшающими поисковые возможности методов локальной оптимизации.

Свойство 3.2. Если $y_{\gamma_k}^* (\varepsilon_k) \in E \setminus Y$, т.е. оценки приближаются к Y^* извне допустимого множества, то любая предельная точка y_∞^* будет находиться на пересечении границ $\partial Y^* \cap \partial Y$.

ДОКАЗАТЕЛЬСТВО [5]. В теореме 3.1 показано, что $y_\infty^* \in Y^* \subseteq Y$. Поэтому в условиях свойства 3.2 найдется подпоследовательность точек $y_{\gamma_k}^* (\varepsilon_k)$, $k=k_s$, не принадлежащих Y и Y^* и сходящуюся к y_∞^* . Это доказывает справедливость свойства.

Свойство 3.3. Если $Y^* \subseteq \text{int } Y$ и штрафные задачи решаются точно (т.е. $\varepsilon_k=0$), то $\exists k$, что $y_{\gamma_k}^* \in Y^*$.

ДОКАЗАТЕЛЬСТВО [5]. По теореме 3.1 для $\delta > 0$ такого, что $O_\delta(Y^*) \subseteq Y \exists k=k(\delta)$, что $\rho(y_{\gamma_k}^*, Y^*) < \delta$, и поэтому $y_{\gamma_k}^* \in Y$. Для таких k $S_{\gamma_k}(y)$ достигает минимума в точках множества Y . Но для $y \in Y$ $S_{\gamma_k}(y) \equiv f(y)$, значит, эта точка будет глобальным минимумом f на Y , т.е. $y_{\gamma_k}^* \in Y^*$.

Свойство 3.4. Если в методе штрафов последовательность коэффициентов штрафа образует неубывающую последовательность $\gamma_{k+1} \geq \gamma_k$, то последовательность значений $f_k = f(y_{\gamma_k}^*)$ будет неубывающей, а последовательность $H_k = H(y_{\gamma_k}^*)$ — не возрастающей:

$$f(y_{\gamma_{k+1}}^*) \geq f(y_{\gamma_k}^*), \quad H(y_{\gamma_{k+1}}^*) \leq H(y_{\gamma_k}^*). \quad (3.13)$$

ДОКАЗАТЕЛЬСТВО. Рассмотрим два значения коэффициента штрафа $\gamma_{k+1} \geq \gamma_k$. Тогда будут справедливы два очевидных неравенства

$$\begin{aligned} f(y_{\gamma_k}^*) + \gamma_k H(y_{\gamma_k}^*) &\leq f(y_{\gamma_{k+1}}^*) + \gamma_k H(y_{\gamma_{k+1}}^*), \\ f(y_{\gamma_{k+1}}^*) + \gamma_{k+1} H(y_{\gamma_{k+1}}^*) &\leq f(y_{\gamma_k}^*) + \gamma_{k+1} H(y_{\gamma_k}^*). \end{aligned}$$

После преобразования каждого из них получим, что

$$\begin{aligned} 0 &\leq (f(y_{\gamma_{k+1}}^*) - f(y_{\gamma_k}^*)) + \gamma_k (H(y_{\gamma_{k+1}}^*) - H(y_{\gamma_k}^*)), \\ (f(y_{\gamma_{k+1}}^*) - f(y_{\gamma_k}^*)) + \gamma_{k+1} (H(y_{\gamma_{k+1}}^*) - H(y_{\gamma_k}^*)) &\leq 0. \end{aligned}$$

Заметим, что полученные выражения совпадают с точностью до значения коэффициента штрафа.

Таким образом, мы видим, что при увеличении коэффициента штрафа со значения γ_k до γ_{k+1} приведенное выше выражение из неотрицательного становится неположительным. Это возможно только в том случае, когда $(f(y_{\gamma_{k+1}}^*) - f(y_{\gamma_k}^*)) \geq 0$ и $(H(y_{\gamma_{k+1}}^*) - H(y_{\gamma_k}^*)) \leq 0$. Тем самым теорема доказана.

Заметим, что в том случае, когда решение исходной задачи размещается на границе допустимой области Y и не является точкой безусловного минимума

целевой функции, из теоремы следует, что решения вспомогательных штрафных задач приближаются к точке решения извне этой области.

Рассмотрим еще один важный аспект — возможность конечной верхней оценки для коэффициента штрафа, при допущении ненулевой малой невязки в ограничениях по значению функции штрафа. Это вопрос исследуется в следующем свойстве.

Свойство 3.5. Пусть известна хотя бы одна допустимая точка $y_Y \in Y$ и $\forall k: \gamma_{k+1} \geq \gamma_k$. Если в качестве критерия останова в методе штрафов принять выполнение условия $H(y_{\gamma_k}^*) \leq \varepsilon_H$ ($\varepsilon_H > 0$), то будет существовать конечная оценка значения коэффициента штрафа γ^* , такая, что при

$$\gamma_k \geq \gamma^* = (f(y_Y) - f(y_{\gamma_0}^*)) / \varepsilon_H \quad (3.14)$$

точное решение задачи со штрафом будет удовлетворять указанному критерию останова.

ДОКАЗАТЕЛЬСТВО. Запишем оценку, используя свойство 3.4

$$H(y_{\gamma_k}^*) = \frac{(S_{\gamma_k}(y_{\gamma_k}^*) - f(y_{\gamma_k}^*))}{\gamma_k} \leq \frac{(f(y_Y) - f(y_{\gamma_k}^*))}{\gamma_k} \leq \frac{f(y_Y) - f(y_{\gamma_0}^*)}{\gamma_k}.$$

Потребуем, чтобы данная оценка не превосходила ε_H . Это будет достигаться $\forall \gamma_k \geq \gamma^*$. Свойство доказано.

3.2.4. Оценки скорости сходимости метода внешнего штрафа

Приведем результаты по скорости сходимости, сделав дополнительные предположения о свойствах задачи.

Дополнительные предположения.

- A. Отсутствует погрешность решения задач со штрафом ($\varepsilon_k \equiv 0$).
- B. Используется степенной штраф вида (3.5) с $r > 0$, а невязка в ограничениях оценивается функцией $G(y)$ из (3.12).
- C. Множество E — замкнуто, а f — липшицева на E с константой L в метрике $\rho(\cdot, \cdot)$.
- D. $\exists \delta > 0$ и $\alpha > 0$, что $\forall y \in (O_\delta(Y) \setminus Y) \cap E$ выполняется $G(y) \geq \alpha \cdot \rho(y, Y)$.

Будем измерять ошибку оценивания решения y^* исходной задачи (3.1), (3.2) через решение y_γ^* вспомогательной задачи (3.7), используя разность значений функций исходной и вспомогательной задач. Запишем величину ошибки

$$\Delta(\gamma) = f(y^*) - S_\gamma(y_\gamma^*) \quad (3.15)$$

Заметим, что по свойству 3.4 всегда $\Delta(\gamma) \geq 0$.

Теорема 3.2 (об оценке скорости сходимости). Если выполнены условия 1–4 теоремы сходимости 3.1 и кроме того — дополнительные предположения A–D, то при точном решении вспомогательных задач со штрафом

- a) при $r \leq 1$ существует достаточно большое γ , что $\forall \gamma \geq \gamma$ выполнится $\Delta(\gamma) = 0$;
- в) при $r > 1$ и достаточно больших γ $0 \leq \Delta(\gamma) \leq C / \gamma^{1/(r-1)}$, где

$$C = (1 - 1/r)(L/\alpha)^{r/(r-1)} / r^{1/(r-1)}. \quad (3.16)$$

ДОКАЗАТЕЛЬСТВО [5]. Следует рассмотреть два случая.

А. Пусть при некотором γ решение задачи со штрафом (3.5)–(3.7) $y_\gamma^* \in Y$ тогда $y_\gamma^* \in Y^*$ (поскольку для $y \in Y$ $S_\gamma(y) \equiv f(y)$). Следовательно $\Delta(\gamma) = 0$.

В. Во всех остальных случаях $\forall \gamma : y_\gamma^* \in E \setminus Y$. Рассмотрим эти случаи. По теореме 3.1 при $\gamma \rightarrow \infty$ $\rho(y_\gamma^*, Y^*) \rightarrow 0$. Тогда $\forall \delta > 0$ при достаточно больших γ : $y_\gamma^* \in O_\delta(Y^*) \setminus Y^*$. По дополнительному предположению (D)

$$G(y_\gamma^*) \geq \alpha \cdot \rho(y_\gamma^*, Y^*).$$

Рассмотрим точку z_γ — проекцию y_γ^* на Y^* :

$$z_\gamma = \pi_{Y^*}(y_\gamma^*).$$

Выполним оценку величины штрафа

$$H(y_\gamma^*) \geq \left(\max \left\{ \max \{ 0; c_1 g_1(y_\gamma^*); \dots; c_m g_m(y_\gamma^*) \}; \max \{ |C_s h_s(y_\gamma^*)| : s = 1, \dots, p \} \right\} \right)^r = (G(y_\gamma^*))^r \geq \alpha^r (\rho(y_\gamma^*, Y^*))^r.$$

Теперь оценим значение минимума функции задачи со штрафом

$$S_\gamma(y_\gamma^*) = f(y_\gamma^*) + \gamma \cdot H(y_\gamma^*) \geq f(z_\gamma) - |f(z_\gamma) - f(y_\gamma^*)| + \gamma \alpha^r (\rho(y_\gamma^*, Y^*))^r \geq f(y^*) - L\rho(y_\gamma^*, Y^*) + \gamma \alpha^r \rho(y_\gamma^*, Y^*)^r = f(y^*) - L\rho + \gamma \alpha^r \rho^r,$$

где $\rho = \rho(y_\gamma^*, Y^*)$.

Следовательно, $0 \leq \Delta(\gamma) \leq \rho(L - \gamma \alpha^r \rho^{r-1})$. Вид зависимости для разных диапазонов значений r показан на рис.3.4.

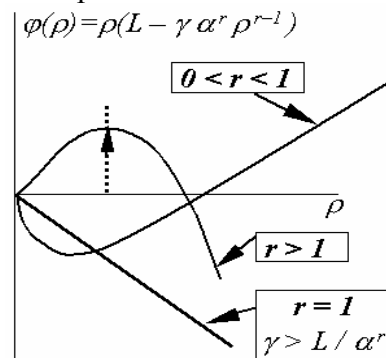


Рис.3.4. Вид верхних оценок для $\Delta(\gamma)$ при различных r

Рассмотрим несколько случаев.

Пусть $0 < r < 1$. При $\gamma \rightarrow \infty$ $\rho = \rho(y_\gamma^*, Y^*) \rightarrow 0$, следовательно $\gamma \alpha^r \rho^{r-1} \rightarrow \infty$. Начиная с некоторого $\bar{\gamma}$ верхняя оценка для $\Delta(\gamma)$ станет отрицательна, что противоречит возможному знаку $\Delta(\gamma)$. Это означает, что начинают нарушаться условия, при которых эта оценка получена, т.е. окажется выполненным $y_\gamma^* \in Y$. Тогда, согласно рассмотренному случаю (А), $y_\gamma^* \in Y^*$, а значит, $\Delta(\gamma) = 0$.

Пусть $r = 1$. При $\gamma > L/\alpha^r$ $\Delta(\gamma) < 0$, если $\rho > 0$. Повторяя приведенные выше рассуждения, приходим к выводу, что $\Delta(\gamma) = 0$.

Пусть $r > 1$. В этом случае нетрудно показать (см. рис. 3.4), что при указанном в (3.16) значении C выполнится неравенство

$$\Delta(\gamma) \leq \max \{ \rho (L - \gamma \alpha^r \rho^{r-1}) : \rho \geq 0 \} = C/\gamma^{1/(r-1)}.$$

Возникает вопрос о том, при каких условиях можно точно гарантировать, что ограничения в задаче удовлетворяют дополнительному условию (D).

Теорема 3.3. Если $g_j(y)$ для $j=1, \dots, m$ выпуклы и непрерывны на выпуклом компакте E , область $\{y \in E: g(y) \leq 0\}$ удовлетворяет условию Слейтера: $\exists \bar{y} \in E$, что $g(\bar{y}) < 0$, для ограничений–равенств $\forall y \in E$, что $h(y) = 0$, выполняется $\nabla h(y) \neq 0$, и h – непрерывно дифференцируема на E , то $\exists \delta > 0, \alpha > 0$, что $\forall y \in O_\delta(Y) \setminus Y \cap E$ невязка $G(y) \geq \alpha \rho(y, Y)$.

ДОКАЗАТЕЛЬСТВО. Очевидно, что если обозначить через $\beta = \min \{ \|\nabla h_s(y)\| : s=1, \dots, p, y \in E, h(y)=0 \}$, то $\beta > 0$. Поэтому при некотором $\delta > 0$ для ограничений–равенств в невязке $G(y)$ из (3.12) следует выбрать, например, $\alpha \leq \min \{ 0,5\beta |C_s| : s=1, \dots, p \}$.

Рассмотрим неравенства. Пусть $G_g(y) = \max\{g_j(y) : j=1, \dots, m\}$. По условию $\exists \bar{y} \in E$, что $G_g(\bar{y}) < -d < 0$. Пусть в точке y из E нарушаются ограничения–неравенства, отсюда $G_g(y) > 0$. Очевидно (см. рис. 3.5), что $\exists y' \in [y, \bar{y}]$, что $G_g(y') = 0$.

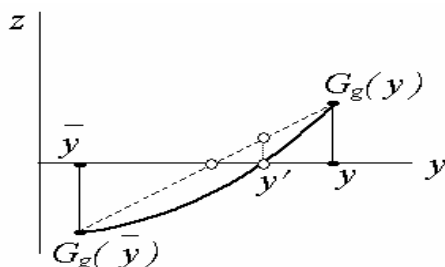


Рис.3.5. Поведение невязки у границ области при выпуклых ограничениях

Из выпуклости: $G_g(y) / \|y - y'\| \geq -G_g(\bar{y}) / \|y' - \bar{y}\|$, поэтому

$$G_g(y) \geq -G_g(\bar{y}) \cdot \|y - y'\| / \|y' - \bar{y}\| \geq (d / \text{diam } E) \rho(y, Y).$$

Отсюда и из ранее сделанной оценки для ограничений–равенств следует справедливость теоремы.

3.2.5. Недостаточность локальных методов при использовании метода штрафов

В теореме об условиях сходимости метода штрафов предполагается, что при решении каждой задачи со штрафом определяется $y_{\gamma_k}^*(\varepsilon_k)$ — оценка глобального минимума штрафной задачи. Поскольку методы многоэкстремальной оптимизации требуют значительно большего объема вычислений, чем методы локальной оптимизации, то при практических расчетах, обычно прибегают к следующему приему. На первой итерации метода штрафов для вычисления $y_{\gamma_0}^*(\varepsilon_0)$ используют один из методов многоэкстремальной оптимизации, а на следующих итерациях для получения оценок $y_{\gamma_{k+1}}^*(\varepsilon_{k+1})$ ($k=1, \dots$) прибегают к методам локальной оптимизации, в которых в качестве начальных точек поиска используются точки $y_{\gamma_k}^*(\varepsilon_k)$, найденные на предыдущей итерации.

Таким образом, полностью отказываться от применения методов глобального поиска нельзя, они необходимы хотя бы на первой итерации метода штрафов. Можно указать две ситуации, в которых ошибка в определении

начальной оценки $y_{\gamma_0}^*(\varepsilon_0)$ может привести к последующей потере решения. Первая соответствует тому случаю, когда решение y^* задачи с ограничениями является внутренней точкой множества Y . При этом, если оценка $y_{\gamma_0}^*(\varepsilon_0)$ не будет принадлежать области притяжения решения y^* , то при последующем использовании локальных методов оценка решения не будет приближаться к точному решению. Похожий эффект возможен и при расположении решения y^* на границе допустимой области, если начальная оценка $y_{\gamma_0}^*(\varepsilon_0)$ окажется в окрестности локального минимума функции штрафа $H(y)$, расположенного вне допустимой области. Действительно, при всех достаточно больших значениях коэффициента штрафа γ у $S_\gamma(y)$ — функции штрафной задачи, будет существовать локальный минимум, расположенный в малой окрестности локального минимума функции штрафа. Если такая ситуация возникла на первой же итерации метода штрафов, и оценка $y_{\gamma_0}^*(\varepsilon_0)$ попала в область притяжения этого локального минимума, то при последующем использовании локальных методов решения задач со штрафом оценка решения останется в окрестности указанного локального минимума и не будет приближаться к точному решению. Характерным признаком, по которому вычислительный метод может распознать эту ситуацию, является неограниченный рост значений функции $S_{\gamma_k}(y)$, вычисляемых в точках получаемых оценок $y_{\gamma_k}^*(\varepsilon_k)$, при $k \rightarrow \infty$.

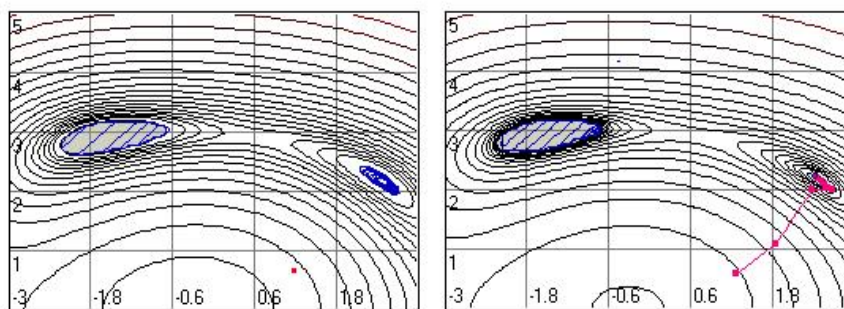


Рис.3. 6 *Возможность потери решения при наличии локального минимума у функции штрафа*

На рис.3.6 приведен пример описанной выше ситуации. В этом примере решается задача поиска минимума линейной функции $-y_1 + y_2$ в прямоугольной области $-3 \leq y_1 \leq 3, 0 \leq y_2 \leq 5$ при дополнительном ограничении $(y_1^2 + y_2^2 - 11)^2 + (y_1 + y_2^2 - 7)^2 + (y_1 - y_2) \leq 0.3$. Функция штрафа в этом примере имеет локальный минимум, расположенный в окрестности точки $y_1 = 2,45, y_2 = 2,15$. При сколь угодно большем увеличении коэффициента штрафа у функции штрафной задачи сохраняется локальный минимум в окрестности этой точки. На рисунке показано поведение метода прямого поиска Хука–Дживса (см. раздел 7.7), запущенного из начальной точки, расположенной в области притяжения этого минимума. Видно, что правильное решение оказалось потерянным.

Допустимая область выделена на рисунке более темным цветом, в ней показаны изолинии линейной минимизируемой функции. Левый рисунок соответствует выбору значения коэффициента штрафа $\gamma = 10$, а правый — значения $\gamma = 1000$.

3.3. Метод модифицированных функций Лагранжа

При знакомстве с этой темой важно понимать ее взаимосвязь с условиями экстремума и теорией двойственности, а также причины, приводящие к введению модификаций функций Лагранжа. Правильные концептуальные представления по этим вопросам весьма важны и позволяют лучше понять сам метод и его свойства.

Прежде чем говорить о методе модифицированных функций Лагранжа следует рассмотреть общую схему метода множителей.

3.3.1. Общая схема метода множителей Лагранжа и ее недостатки

Из главы 2 известно, что для «хорошо» устроенных задач математического программирования (3.1) – (3.2), например, для случая выпуклого E , выпуклых на E функций g и f , а также аффинных функций h , при регулярности допустимой области Y критерием глобальной оптимальности точки y^* является (см. теоремы 2.1, 2.2, 2.4) существование седловой точки (y^*, λ^*, μ^*) в области $E \times (R^+)^m \times R^p$ у функции Лагранжа

$$L(y, \lambda, \mu) = f(y) + (\lambda, g(y)) + (\mu, h(y)),$$

откуда

$$L(y^*, \lambda^*, \mu^*) = \max_{(\lambda, \mu) \in (R^+)^m \times R^p} \left(\min_{y \in E} L(y, \lambda, \mu) \right). \quad (3.17)$$

Причем задача на максимум в (3.17) соответствует двойственной задаче математического программирования (2.42) – (2.43).

Из (3.17) следует, что задача оптимизации (3.1)–(3.2) в пространстве размерности N при функциональных ограничениях может быть (при оговоренных выше условиях) сведена к задаче взятия максимина в области простой структуры (без функциональных ограничений) в пространстве размерности $N+m+p$.

Это важное наблюдение может служить основой для разработки вычислительных алгоритмов. При этом первостепенную важность приобретает вопрос о способах поиска решения максиминных задач (3.17). Есть два варианта реализации вычислительных методов.

Первый вариант реализации, использует градиентные направления в пространстве переменных y и (λ, μ) , и фиксированный коэффициент длины шага. Он был предложен в 60-х годах авторами Эрроу и Гурвицем. Их итерационный метод поиска седловой точки функции Лагранжа имеет вид

$$y^{k+1} = \pi_E(y^k - x \nabla_y L(y^k, \lambda^k, \mu^k)) \quad (3.18)$$

$$\lambda^{k+1} = (\lambda^k + x g(y^k))_+ \quad (3.19)$$

$$\mu^{k+1} = \mu^k + x h(y^k). \quad (3.20)$$

Здесь $g(y^k)$ и $h(y^k)$ являются векторами градиентов функции Лагранжа по переменным λ и μ и, одновременно по теореме (2.14), — компонентами субградиента двойственной функции Лагранжа (2.42). Операторы π_E и $(\cdot)_+$ выполняют проектирование на множества E и $(R^+)^m$. Значение коэффициента шага $x > 0$ является параметром метода.

К сожалению, процедура (3.18)–(3.20) может порождать расходящийся процесс. Чтобы показать это достаточно привести пример.

ПРИМЕР РАСХОДИМОСТИ. Рассмотрим простую задачу с $E=R^1$, $m=1$, $g_1(y) = -y \leq 0$, $f(y)=y$, ограничения–равенства отсутствуют (т.е. $p=0$). Следует заметить, что эта задача удовлетворяет всем указанным выше условиям. Ее решением является точка $y^*=0$, а соответствующее значение λ^* в условиях экстремума равно 1.

Поскольку функция Лагранжа в этой задаче имеет вид $L(y, \lambda) = y - \lambda y$, то без учета операции проектирования $(\cdot)_+$ итерационный процесс (3.18)–(3.20) примет для нее следующую форму

$$y^{k+1} = y^k - x(1 - \lambda^k), \quad \lambda^{k+1} = \lambda^k - x y^k. \quad (3.21)$$

Выразив y^k из второго уравнения и исключив за счет этого y^{k+1} , y^k в первом, найдем следующую связь между последовательными значениями λ :

$$\lambda^{k+1} - \lambda^{k+2} = \lambda^k - \lambda^{k+1} - x^2(1 - \lambda^k)$$

или

$$\lambda^{k+2} - 2\lambda^{k+1} + \lambda^k(1 + x^2) = x^2 \quad (3.22)$$

Стационарным решением последнего разностного уравнения является значение $\lambda^* = 1$, что правильно соответствует значению множителя Лагранжа в условиях экстремума. Поскольку $\lambda^* > 0$, то при начальном значении λ^0 близком к 1 опущенный оператор проектирования $(\cdot)_+$ для λ не будет оказывать влияния на итерационный процесс, пока он выполняется в области $\lambda > 0$. Поэтому сходимость процесса можно исследовать без его учета.

Необходимым и достаточным условием асимптотической устойчивости стационарного решения $\lambda^k \equiv \lambda^*$ для уравнения (3.22) является размещение внутри единичного круга корней характеристического полинома $z^2 - 2z + (1 - x^2) = 0$, но $z_{1,2} = 1 \pm (1 - (1 - x^2))^{0.5} = 1 \pm |x|$, следовательно $\max |z_{1,2}| = 1 + |x| \geq 1$. Процесс итераций расходится для любого x . Возможность расходимости (3.18)–(3.20) доказана.

Второй вариант реализации метода множителей непосредственно основан на решении двойственной задачи (2.42)–(2.43) с использованием той же итерационной схемы по λ и μ , что и в методе Эрроу–Гурвица, но с полным вычислением двойственной функции Лагранжа. Его итерационные формулы приведены в конце главы 2 в формулах (2.50)–(2.52).

Необходимо сразу подчеркнуть, что область его применения ограничивается задачами, где нет разрыва двойственности (см. теорему 2.11 и рис. 2.5). Кроме того, этот вычислительный метод нельзя применить при неограниченных значениях двойственной функции Лагранжа. Например, в рассмотренном выше примере $L^*(\lambda, \mu) = -\infty$ при $\lambda \neq 1$ и 0 – при $\lambda = 1$. Хотя это и не противоречит (3.17), но препятствует поиску решения численными методами.

Таким образом, простой метод множителей в общем случае не может быть применен, он нуждается в изменениях. Для того, чтобы понять пути его модификации, полезно видоизменить форму постановки задачи так, чтобы метод было проще анализировать.

3.3.2. Преобразование постановки задачи, сведение задач с неравенствами к задачам с равенствами

В целях удобства дальнейшего изложения будем считать, что множество $E=R^N$, т.е. задача ставится так:

$$f(y) \rightarrow \min, y \in Y = \{y \in \mathbb{R}^N : g(y) \leq 0; h(y) = 0\} \quad (3.23)$$

Тем самым двусторонние ограничения на переменные, если они есть, переводятся в разряд функциональных ограничений.

Задачу будем считать достаточно гладкой, что позволит записывать условия экстремума в форме градиентных равенств. Наличие дополнительных требований вида $y \in E \subset \mathbb{R}^N$ лишило бы нас возможности использовать эту форму записи в общем случае.

Чтобы упростить изложение, избавимся от ограничений–неравенств. Существует общий прием сведения неравенств к равенствам за счет увеличения размерности пространства переменных. Рассмотрим неравенство $g_i(y) \leq 0$. Введем дополнительные переменные z_i , чтобы $g_i(y) = -z_i^2$. Отсюда эквивалентной формой записи неравенства будет

$$g_i(y) + z_i^2 = 0.$$

Будем считать, что пространство расширено за счет добавления к y переменных z , и выполним замены

$$y := (y \mid \underbrace{z_1, \dots, z_m}_z) \in \mathbb{R}^{N+m}, \quad N := N + m.$$

Окончательно приходим к задачам, в которых присутствуют только ограничения–равенства

$$f(y) \rightarrow \min, y \in Y = \{y \in \mathbb{R}^N : h(y) = 0\}. \quad (3.24)$$

Функция Лагранжа будет иметь вид

$$L(y, \mu) = f(y) + (\mu, h(y)) \quad (3.25)$$

3.3.3. Построение модифицированной функции Лагранжа

Вначале вернемся к общей постановке задачи (3.1)–(3.2) (позднее результаты проведенного анализа будут применены к постановке (3.24)).

Целью модификации функции Лагранжа является расширение области применения метода множителей, а также улучшение сходимости возникающих итерационных схем. В первую очередь будем исходить из цели расширения области применимости.

Заметим, что условие существования седловой точки $y \in L(y, \lambda, \mu)$ не является необходимым условием для задач общего вида, а также для произвольных задач вида (3.24). Таким образом, в точке решения y^* у функции $L(y, \lambda^*, \mu^*)$ может не быть минимума по y на E .


Пример. Рассмотрим задачу с невыпуклой функцией $f(y) = y^3$ и одним ограничением $h(y) = y + 1 = 0$:

$$y^3 \rightarrow \min, y \in Y = \{y \in \mathbb{R}^1 : y + 1 = 0\}.$$

В этой задаче существует единственная допустимая точка $y = -1$. Условие экстремума дает $-3(-1)^2 = \mu^* \cdot 1$, следовательно $\mu^* = -3$, $L(y, \mu^*) = y^3 - 3(y + 1)$.

В точке $y = y^* = -1$ $\nabla_y L(y^*, \mu^*) = 0$, $\Gamma_y^L(y^*, \mu^*) = -6 < 0$, следовательно в этой точке по переменной y у функции Лагранжа достигается максимум, а не минимум.

Таким образом, единственное, что можно гарантировать у функций Лагранжа $L(y, \lambda, \mu)$ для гладких регулярных задач в точке решения $y^* \in E$ – это стационарность точки y^*, λ^*, μ^* .

 **Поставим задачу:** так модифицировать функцию Лагранжа, чтобы точки (y^*, λ^*, μ^*) оставались седловыми, но при $\lambda = \lambda^*$, $\mu = \mu^*$ модифицированная функция имела хотя бы локальный минимум по y в точке y^* , даже если задача не выпукла.

Рассмотрим структуру функции Лагранжа. Пусть для y^* известен набор активных ограничений $J^* = J(y^*)$ и значения самих множителей Лагранжа λ^*, μ^* . В силу условий Куна–Таккера $\lambda_j^* g_j(y^*) = 0$ ($j = 1, \dots, m$) и, кроме того, $h(y^*) = 0$. Поэтому на множестве

$$Y(y^*) = \{y \in E: h(y) = 0; J(y) = J(y^*); g_i(y) < 0, i \notin J(y)\}$$

$$L(y, \lambda^*, \mu^*) = f(y), \quad y \in Y(y^*). \quad (3.26)$$

Поскольку y^* — точка минимума функции f на Y , и $y^* \in Y(y^*) \subseteq Y$, то на множестве $Y(y^*)$ функция Лагранжа (в силу (3.26)) также будет иметь в точке y^* минимум.

Для того, чтобы данное свойство сохранилось не только на $Y(y^*)$, но и на всем множестве E (или хотя бы локально, в области $E \cap O_\varepsilon(y^*)$) необходимо сделать добавку к $L(\cdot)$, обеспечивающую возрастание этой функции при удалении от многообразия $Y(y^*)$.

Чтобы упростить изложение, вернемся к постановке (3.24), т.е. будем считать, что задача представлена в форме задачи с равенствами. Наложим на нее дополнительные ограничения, гарантирующие существование в точке y^* строгого локального минимума f на Y . Для этого воспользуемся достаточными условиями второго порядка для строго локального минимума из теоремы 2.8.

Предположения о задаче (3.24)

A. Пусть $f, h \in C^2(R^N)$, y^* — решение задачи (3.24), являющееся ее регулярной точкой.

B. Для множителей Лагранжа μ^* , при которых выполнено условие экстремума первого порядка $\nabla_y L(y^*, \mu^*) = 0$, дополнительно имеет место следующее:

$$\forall d \neq 0, \text{ что } (d, \nabla h_s(y^*)) = 0 \quad (s = 1, \dots, p): d^T \Gamma_y^L(y^*, \mu^*) d > 0. \quad (3.27)$$

При этих предположениях y^* будет являться строгим локальным минимумом в задаче (3.24). Условие (3.27) можно трактовать как условие положительной определенности матрицы Гессе по переменным y для функции Лагранжа $L(y, \mu^*)$ на линейном многообразии касательных направлений

$$K^H(y^*) = \{d \in R^N : (d, \nabla h_s(y^*)) = 0 \quad (s = 1, \dots, p)\},$$

которое было введено при доказательстве теоремы 2.7 (см. формулу (2.28)).

Заметим, что модифицированная функция Лагранжа должна иметь в точке y^* матрицу Гессе (гессиан), положительно определенную на всем R^N .

Построим модифицированную функцию Лагранжа $L_\gamma(y, \mu)$, добавив к $L(\cdot)$ слагаемое, возрастающее при уходе точки y с ограничений–равенств:

$$L_\gamma(y, \mu) = f(y) + (\mu, h(y)) + \frac{1}{2} \gamma \|h(y)\|^2 \quad (3.28)$$

ИССЛЕДОВАНИЕ. Изучим эту функцию.

$$\nabla_y L_\gamma(y, \mu) = \nabla f(y) + (\mu + \gamma h(y))^T \nabla h(y)$$

$$\begin{aligned} \Gamma_y^{L_\gamma}(y^*, \mu^*) &= \Gamma^f(y^*) + \sum_{j=1}^p (\mu_j^* + \underbrace{\gamma h_j(y^*)}_{=0}) \Gamma^{h_j}(y^*) + \gamma (\nabla h(y^*))^T \nabla h(y^*) = \\ &= \Gamma_y^L(y^*, \mu^*) + \gamma (\nabla h(y^*))^T \nabla h(y^*). \end{aligned}$$

Введем два обозначения:

$$V(d) = d^T \Gamma_y^L(y^*, \mu^*) d, \quad W(d) = d^T (\nabla h(y^*))^T \nabla h(y^*) d = \|\nabla h(y^*) d\|^2.$$

Тогда

$$\forall d \neq 0, d \in K^H(y^*): \quad V(d) > 0, W(d) = 0$$

$$\forall d \neq 0, d \notin K^H(y^*), \quad W(d) > 0.$$

Для достаточно малого ε и $\forall \xi \in O_\varepsilon(0)$ при $d \in K^H(y^*)$ выполнится $V(d+\xi) > 0$, $W(d+\xi) \geq 0$. Обозначим

$$\eta_1 = \inf \{ V(d): \|d\|=1, d \in (K^H(y^*) + O_\varepsilon(0)) \},$$


$$\eta_2 = \inf \{ W(d): \|d\|=1, d \notin (K^H(y^*) + O_\varepsilon(0)) \}.$$


Заметим, что множество $K^H(y^*) + O_\varepsilon(0)$ представляет ε -слой вокруг линейного многообразия $K^H(y^*)$. Поскольку $\eta_2 > 0$, то $\exists \gamma > 0$, что $\eta_1 + \gamma \eta_2 > 0$.

Все это, в совокупности, определяет положительную определенность гессиана функции $L_\gamma(y, \mu^*)$ в точке y^* . Замечая, что $\nabla_y L_\gamma(y^*, \mu^*) = \nabla_y L(y^*, \mu^*) = 0$, видим, что для $L_\gamma(y, \mu^*)$ выполнены достаточные условия строгого локального минимума в точке y^* по переменным y .

Лемма 3.4. Для задачи (3.24), удовлетворяющей предположениям A, B при достаточно большом γ модифицированная, согласно (3.28), функция Лагранжа имеет в точке y^* строгий локальный минимум по y .

Доказательство приведено перед текстом леммы.

 **Замечание 1.** В (3.28) модификация функции Лагранжа выполняется введением штрафной добавки, подобно тому, как это было сделано в методе штрафных функций.

 **Замечание 2.** Формула (3.28) представляет лишь один из возможных способов модификации функции Лагранжа. В литературе по этому вопросу можно встретить гораздо более общие по сравнению с (3.28) выражения, также называемые модифицированными функциями Лагранжа (см., например, [5,23]).

3.3.4. Метод модифицированной функции Лагранжа для задач с ограничениями–равенствами

Метод модифицированной функции Лагранжа имеет следующую итерационную форму

$$y^{k+1} = \arg \min \{ L_{\gamma_k}(y, \mu^k): y \in R^N \} \quad (3.29)$$

$$\mu^{k+1} = \mu^k + \gamma_k h(y^{k+1})$$

$$\gamma_{k+1} \geq \gamma_k \quad (3.30)$$

В отличие от метода штрафов, последовательность γ_k не обязана стремиться к бесконечности. Сходимость может быть обеспечена за счет настройки множителей μ^k .

Теорема 3.5. Пусть для задачи с ограничениями–равенствами (3.24) выполнены предположения A, B, тогда $\exists \bar{\gamma}$ и $\delta > 0$, что при выполнении условий

$$|\mu_0 - \mu^*| < \delta \gamma_0 \quad \text{и} \quad \bar{\gamma} \leq \gamma_k \leq \gamma_{k+1} \quad \forall k$$

для метода модифицированной функции Лагранжа (3.29)–(3.30):

1) вспомогательные задачи (3.29) будут иметь единственное решение, причем $\mu^k \rightarrow \mu^*$ и $y^k \rightarrow y^*$ при $k \rightarrow \infty$;

2) если $\lim_{k \rightarrow \infty} \gamma_k = \gamma^* < \infty, \forall k : \mu^k \neq \mu^*$, то $\exists 0 \leq q < 1$, что

$$\limsup_{k \rightarrow \infty} \left| \frac{\mu^{k+1} - \mu^*}{\mu^k - \mu^*} \right| \leq q \quad (\text{линейная сходимость})$$

3) если $\gamma_k \rightarrow \infty$ при $k \rightarrow \infty$ и $\forall k \mu^k \neq \mu^*$, то

$$\lim_{k \rightarrow \infty} \left| \frac{\mu^{k+1} - \mu^*}{\mu^k - \mu^*} \right| = 0 \quad (\text{сверхлинейная сходимость}).$$

ДОКАЗАТЕЛЬСТВО теоремы требует обоснования ряда вспомогательных фактов и из-за ограничений объема не приводится. Его можно найти в книге Д.Бертсекас [2], раздел 2.2.

Следует обратить внимание на сохранение сходимости при ограниченных значениях коэффициента штрафа γ . Это свойство метода следует считать положительным, поскольку именно неограниченный рост коэффициента штрафа является основным фактором, ухудшающим структуру вспомогательных задач в методе штрафных функций. Однако метод (3.29)–(3.30) имеет свои отрицательные стороны, а именно, его сходимость гарантируется только при достаточно точном начальном приближении множителей Лагранжа.

Применение метода требует автоматической настройки коэффициента штрафа. Здесь может быть применен подход, похожий на используемый в методе штрафов, а именно.

АЛГОРИТМ НАСТРОЙКИ КОЭФФИЦИЕНТА ШТРАФА. Введем параметр $\alpha < 1$ — коэффициент желаемого уменьшения невязки в ограничениях на шаге метода, а также множитель $\beta > 1$ прироста коэффициента штрафа, например, $\alpha=0.25, \beta=10$. Обычно используется следующая схема настройки γ_k

$$\gamma_{k+1} = \begin{cases} \beta \gamma_k, & \|h(y^{k+1})\| \geq \alpha \|h(y^k)\| \\ \gamma_k, & \|h(y^{k+1})\| < \alpha \|h(y^k)\|. \end{cases} \quad (3.31)$$

Таким образом, при достаточно быстром убывании невязки коэффициент штрафа расти не будет.

3.3.5. Метод модифицированной функции Лагранжа в задачах с равенствами и неравенствами

Вернемся от задачи (3.24) к задаче с равенствами и неравенствами (3.23), но воспользуемся в явном виде рассмотренным ранее приемом сведения ограничений–неравенств к форме равенств $g_i(y) + z_i^2 = 0 \quad (i=1, \dots, m)$.

Построим модифицированную функцию Лагранжа в виде

$$L_\gamma(y, z, \lambda, \mu) = f(y) + (\mu, h(y)) + 0.5 \gamma \|h(y)\|^2 + \sum_{i=1}^m \{ \lambda_i (g_i(y) + z_i^2) + 0.5 \gamma (g_i(y) + z_i^2)^2 \}. \quad (3.32)$$

Минимум в (3.29) по переменным (y, z) разделяется в последовательное взятие минимумов по z и по y . Находя минимум по z в (3.32) аналитически, получим

$$(z_i^*)^2 = \max\{0; -\lambda_i/\gamma - g_i(y)\}$$

Подставив это выражение в (3.32), получим окончательный вид модифицированной функции Лагранжа с исключенной переменной z :

$$L_\gamma(y, \lambda, \mu) = f(y) + (\mu, h(y)) + 0.5 \gamma \|h(y)\|^2 + \sum_{i=1}^m \left(\lambda_i \max\{g_i(y); -\frac{\lambda_i}{\gamma}\} + 0.5 \gamma \left(\max\{g_i(y); -\frac{\lambda_i}{\gamma}\} \right)^2 \right). \quad (3.33)$$

Вычислительная схема метода приобретает следующий вид

$$y^{k+1} = \arg \min\{L_{\gamma_k}(y, \lambda^k, \mu^k) : y \in R^N\}, \quad (3.34)$$

$$\mu^{k+1} = \mu^k + \gamma_k h(y^{k+1}), \quad (3.35)$$

$$\lambda_i^{k+1} = (\lambda_i^k + \gamma_k g_i(y^{k+1}))_+ \quad (3.36)$$

$$\gamma_{k+1} \geq \gamma_k.$$

Настройка коэффициента штрафа выполняется аналогично (3.31), но в невязке по ограничениям необходимо также учитывать меру нарушения неравенств.

3.4. Другие общие методы учета ограничений

В этом разделе будет описана группа методов, которые можно было бы назвать методами параметризации, хотя это и не является общепринятым по отношению к некоторым из них.

Вначале рассмотрим классические методы, обычно относимые к этой группе, а затем — менее известные методы, используемые, обычно в сочетании со специальными алгоритмами глобального поиска, которые будут разобраны позднее, в главах 4–5. Все методы этого пункта изложим применительно к задаче, включающей только ограничения–неравенства

$$f(y) \rightarrow \min, y \in Y, f: Y \rightarrow R^1, \quad (3.37)$$

$$Y = \{y \in E \subseteq R^N : g(y) \leq 0\}. \quad (3.38)$$

Как и в методе штрафов предполагается, что множество E имеет простую структуру, позволяющую при отсутствии функциональных ограничений отыскивать численными методами минимумы функций на E .

3.4.1. Метод параметризации целевой функции

Как и в методе внешнего штрафа здесь вводится штрафная функция $H(y)$ вида (3.4), которая может быть построена аналогично (3.5) в виде степенной функции

$$H(y) = \sum_{i=1}^m ((c_i g_i(y))_+)^r \quad (3.39)$$

с коэффициентами $c_1, \dots, c_m > 0$, $r \geq 0$ (обычно $r = 2$), $(\cdot)_+$ — оператор проектирования на $(\mathbb{R}^+)^1$.

Целевые функции $S_{\beta_k}(y)$ вспомогательных задач (3.3) в методах параметризации строятся так, чтобы параметр β_k влиял не на величину штрафной добавки, а на способ учета целевой функции $f(y)$ исходной задачи. Выбор нового значения параметра β_{k+1} жестко связан в этих методах с полученным оптимальным значением $S_{\beta_k}^*$ функции текущей вспомогательной задачи.

ОПИСАНИЕ МЕТОДА. В качестве начального значения β_0 принимается нижняя оценка глобального минимума функции $f(y)$ на Y , т.е. выбирается $\beta_0 \leq f(y^*)$. Для определения β_0 приходится решать дополнительную экстремальную задачу, например,

$$\beta_0 = \min \{ f(y) : y \in E \}.$$

Далее используются следующие итерационные соотношения

$$S_{\beta_k}(y) = (f(y) - \beta_k)^2 + H(y) \quad (3.40)$$

$$y^k = \arg \min \{ S_{\beta_k}(y) : y \in E \}, \quad (3.41)$$

$$\beta_{k+1} = \Phi(\beta_k, S_{\beta_k}(y^k)) \quad (3.42)$$

где Φ — функция настройки параметра β , может иметь различный вид. Например, при $r = 2$ используется правило

$$\beta_{k+1} = \beta_k + (S_{\beta_k}(y^k))^{1/r} \quad (3.43)$$

На рис. 3.7 приведена геометрическая интерпретация процесса решения задачи (3.37)–(3.38) по методу (3.40), (3.41), (3.43) при двух конфигурациях множеств вида

$$Q_+ = \{ ((c_1 g_1(y))_+, \dots, (c_m g_m(y))_+, f(y)) : y \in E \}.$$

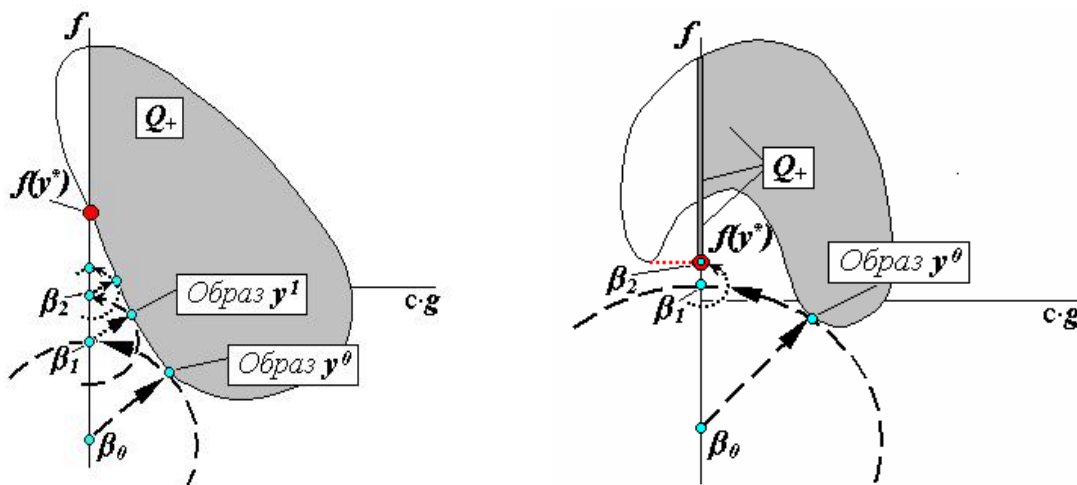


Рис.3. 7. Два типа поведения метода параметризации

Решение задачи (3.41) соответствует отысканию сферы с центром в точке $(0; \beta_k)$ минимального радиуса из имеющих пересечение с множеством Q_+ . Точки

касания этой сферы с Q_+ является образами точек y^k при отображении $((c \cdot g)_+; f): E \rightarrow Q_+$.

Центр новой сферы размещается в точке $(0; \beta_{k+1})$, которая получается в результате пересечения предыдущей сферы с осью f . На рис. 3.7 проиллюстрирован характер процесса сходимости последовательности β_k к оптимальному значению $f^* = f(y^*)$. Левый рисунок соответствует сходимости за бесконечное число шагов, а правый — за конечное. Формальное обоснование метода и условия сходимости определяются следующей теоремой.

Теорема 3.6. Пусть множество $E \subseteq R^N$ является компактом, $f, g \in C(E)$, т.е. непрерывны на E , $Y \neq \emptyset$, $\beta_0 < f(y^*)$. Тогда любая предельная точка y^∞ последовательности y^k ($k=0, 1, 2, \dots$) в методе параметризации (3.40), (3.41), (3.43) при $r=2$ является глобальным минимумом исходной задачи (3.37)–(3.38). Кроме того, существует предел

$$\lim_{k \rightarrow \infty} f(y^k) = f^* = f(y^*). \quad (3.44)$$

ДОКАЗАТЕЛЬСТВО. В силу непрерывности функций и компактности E из теоремы Вейерштрасса вытекает существование решений вспомогательных задач (3.41). Таким образом, последовательность β_k определена. Из (3.43) следует ее монотонность: $\beta_{k+1} \geq \beta_k$. Кроме того, по построению $\beta_0 \leq f^*$.

По индукции легко доказать, что f^* будет являться верхней оценкой для β_k . Действительно, пусть для некоторого k : $\beta_k \leq f^*$. Поскольку $y^* \in Y \subseteq E$, то

$$S_{\beta_k}(y^k) \leq S_{\beta_k}(y^*) \stackrel{\substack{= \\ (c \cdot g(y^*))_+ = 0}}{=} (f^* - \beta_k)^2.$$

Следовательно, $(S_{\beta_k}(y^k))^{1/2} \leq |f^* - \beta_k| = f^* - \beta_k$. Отсюда и из (3.43) вытекает следующая оценка

$$\beta_{k+1} = \beta_k + (S_{\beta_k}(y^k))^{1/2} \leq f^*.$$

Таким образом, β_k не убывает и ограничена сверху значением f^* , следовательно, имеет некоторый предел $\beta^* \leq f^*$. При этом из (3.43) следует, что $S_{\beta_k}(y^k) \rightarrow 0$ при $k \rightarrow \infty$.

Поскольку $S_{\beta_k}(y)$ является суммой неотрицательных слагаемых, то из (3.40) получаем, что

$$f(y^k) \rightarrow \beta^*, \quad H(y^k) \rightarrow 0 \text{ при } k \rightarrow \infty. \quad (3.45)$$

Т.к. $y^k \in E$ и E является компактом, то у этой последовательности существуют предельные точки. Пусть y^∞ — одна из них. Из ранее доказанного следует, что

$$f(y^\infty) = \beta^*, \quad H(y^\infty) = 0.$$

Последнее равенство означает, что точка $y^\infty \in Y$, т.е. допустима. Тогда $f(y^\infty) \geq f^*$, но ранее было показано, что $f^* \geq \beta^*$. Итак: $\beta^* = f(y^\infty) \geq f^* \geq \beta^*$

Отсюда видим, что $f(y^\infty) = f^*$. Теорема доказана.

Замечания. А. Если при некотором k окажется, что $S_{\beta_k}(y^k) = 0$, то найденное решение y^k вспомогательной задачи будет являться точным решением исходной задачи, т.е. $y^k \in Y^*$. При этом поиск следует остановить.

В. Если, за счет вычислительных ошибок, на некотором шаге окажется $\beta_k > f^*$, то $(0, \beta_k) \in Q_+$ и элементы последовательности y^k перестанут изменяться, следовательно, метод потеряет сходимость. Это накладывает повышенные требования на точность решения вспомогательных задач в методе параметризации.

3.4.2. Метод допустимой точки

Метод состоит в следующем. Строится функция

$$G(y) = \max \{ c_1 g_1(y); \dots; c_m g_m(y) \}, \quad (3.46)$$

определяющая верхнюю огибающую системы нормированных (с помощью коэффициентов c_i) функций ограничений.

Замечание. Перед началом использования метода должна быть известна допустимая начальная точка $y^0 \in Y$.

ОПИСАНИЕ МЕТОДА. Вводится параметр метода γ_k . Полагается $\gamma_0 = f(y^0)$. Далее для произвольных $k \geq 0$ определяется итерационная процедура следующего вида

$$S_{\gamma_k}(y) = \max \{ f(y) - \gamma_k; G(y) \} \quad (3.47)$$

$$y^{k+1} = \arg \min \{ S_{\gamma_k}(y) : y \in E \} \quad (3.48)$$

$$\gamma_{k+1} = f(y^{k+1}) \quad (3.49)$$

Заметим, что вспомогательная функция (3.47) напоминает свертку Гермейера для многокритериальных задач (см. пункт 1.3.7, формула (1.31)). Вид функции $S_{\gamma_k}(y)$ для случая $y \in R^1$ для одной из возможных ситуаций показан на рис.3.8 утолщенной линией.

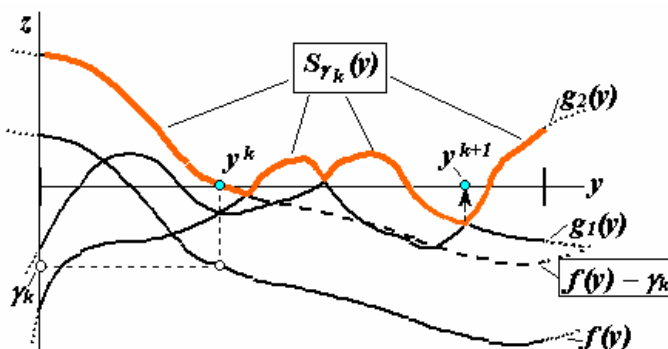


Рис.3. 8. Выполнение итераций в методе допустимой точки

Теорема 3.7. Пусть E — компакт, $f, g \in C(E)$, $Y \neq \emptyset$ и, кроме того, в сколь угодно малой окрестности любой допустимой точки существуют $y \in Y$ со значениями $G(y) < 0$. Тогда любая предельная точка y^∞ последовательности y^k , построенной методом допустимой точки (3.47)–(3.49), является глобальным минимумом задачи (3.37)–(3.38). Кроме того, существует предел, как в (3.44): $\lim_{k \rightarrow \infty} f(y^k) = f^* = f(y^*)$.

ДОКАЗАТЕЛЬСТВО. В силу непрерывности функции (3.47) и компактности E , задачи (3.48) имеют решение.

Легко видеть, что последовательность γ_k является, по построению, не возрастающей, а поскольку все точки y^k допустимы, то $\gamma_k \geq f^* = f(y^*)$, т.е. данная последовательность ограничена снизу. Следовательно, существует предел $\gamma_k \rightarrow \gamma^*$ при $k \rightarrow \infty$. Осталось показать, что $\gamma^* = f(y^*)$.

Пусть это не так и $\gamma^* > f(y^*)$. Рассмотрим значения γ_k достаточно близкие к γ^* :

$\gamma_k - \gamma^* < \delta, \forall k > K$. Далее обратим внимание, что для этих k возможна оценка

$$S_{\gamma_k}(y^{k+1}) \geq f(y^{k+1}) - \gamma_k = \gamma^{k+1} - \gamma^k > -\delta.$$

Рассмотрим точки $\bar{y} \in E$, лежащие в достаточно малой ε -окрестности y^* . По предположению теоремы среди них найдется такая точка \bar{y} что $G(\bar{y}) < -\eta$ для достаточно малого $\eta > 0$. Более того, можно выбрать η и δ настолько малыми, что $\forall k > K$ одновременно выполнится еще одно неравенство: $f(\bar{y}) - \gamma^k < -\eta$. Это возможно в силу непрерывности $f(y)$ и предположения $\gamma^* > f(y^*)$. Таким образом, $S_{\gamma_k}(\bar{y}) < -\eta$.

В силу того, что δ выбирается произвольно и может быть сделано $\delta < \eta$, получим $S_{\gamma_k}(y^{k+1}) > S_{\gamma_k}(\bar{y})$ для $\bar{y} \in E$, что противоречит правилу (3.48). Теорема доказана.

Возможна геометрическая интерпретация метода в пространстве значений функций $(c \cdot g), f$ при $m=1$, показанная на рис.3.9 для двух случаев конфигураций множеств

$$Q = \{ (c_1 g_1(y); \dots; c_m g_m(y); f(y)) : y \in E \}.$$

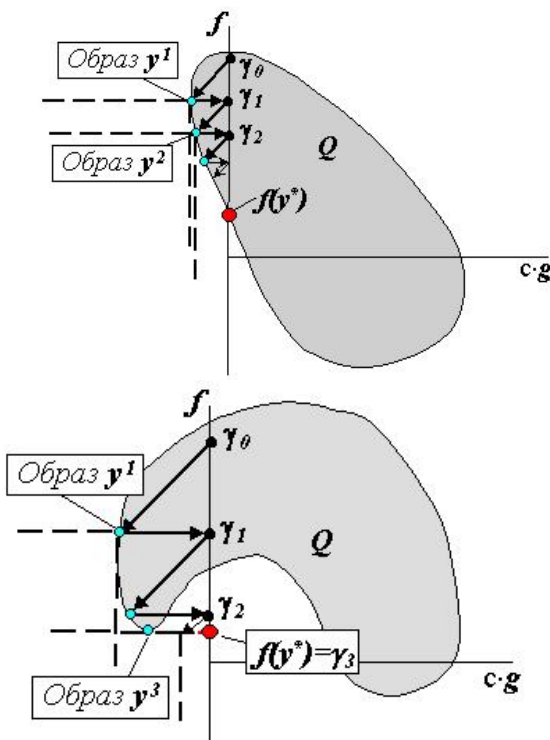



Рис.3. 9. Два типа поведения итераций в методе допустимой точки

Ситуация, представленная на правом рисунке соответствует случаю сходимости за конечное число итераций, на левом — за бесконечное.

 **Замечания.** А. Функции $S_\gamma(y)$ вспомогательных задач метода допустимой точки не являются гладкими, что создает существенные препятствия при использовании методов локальной оптимизации для решения вспомогательных задач (3.48). Этот метод учета ограничений используется, в основном, в специальных алгоритмах многоэкстремальной оптимизации (см. раздел 4.3), где гладкость не является существенной. В. В методах многоэкстремальной оптимизации метод центров обычно используется в модифицированной форме, когда параметр γ_k функции $S_{\gamma_k}(y)$ изменяется непосредственно в процессе поиска ее минимума на E :

$$S_{\gamma_k}(y) \rightarrow \min, y \in E, \quad (3.50)$$

$$\gamma_k = \min \{ f(y^s) : g(y^s) \leq 0, (s=0, \dots, k) \}, \quad (3.51)$$

где y^s — точки выполненных измерений ($s=0, \dots, k$).

3.4.3. Индексный метод учета ограничений

Индексный метод [44] был разработан для одномерных многоэкстремальных задач с частично вычислимыми ограничениями. Его можно рассматривать как развитие метода допустимой точки, описанного в предыдущем пункте, хотя он был разработан независимо. В некотором смысле индексная схема идет дальше метода допустимой точки на пути построения вспомогательных задач, а именно, при этом подходе функции вспомогательных задач оказывается не только не гладкими, но могут оказаться разрывными.

В этом пункте метод будет описан не в своем оригинальном варианте [44], а как общая вычислительная схема учета ограничений.

Метод рассматривается при следующем видоизменении постановки исходной задачи. Обозначим $Y_0 \equiv E$ и определим набор вложенных множеств

$$Y_i = \{ y \in Y_{i-1} : g(y)_i \leq 0 \} \quad (i = 1, \dots, m) \quad (3.52)$$

При этом предполагается, что областью определения i -го ограничения является множество Y_{i-1} , а областью определения целевой функции f — множество Y_m . Введем еще одно вспомогательное пустое множество $Y_{m+1} = \emptyset$, а также введем переобозначение для целевой функции, положив $g_{m+1}(y) \equiv f(y)$. Выполняется следующее включение

$$Y_0 \supseteq Y_1 \supseteq \dots \supseteq Y_m \supseteq Y_{m+1} = \emptyset.$$

Таким образом, часть ограничений и функция f могут быть определены не на всей области E , а лишь на ее части.

Задача оптимизации запишется в виде

$$f(y) \rightarrow \min, y \in Y_m. \quad (3.53)$$

Сопоставим задаче (3.53) вспомогательную задачу, которая будет иметь решение даже в том случае, когда множество Y_m может оказаться пустым. Для этого введем классификацию точек $y \in E$ по значениям индексов $v = v(y)$, определяемых номером первого нарушенного ограничения, точнее $y \in Y_{v-1}$, $y \notin Y_v$. Таким образом, если все ограничения в точке y выполнены, то ее индекс будет равен $m+1$. Запишем это в форме определения.

Определение 3.2. Индексом $\nu(y)$ точки $y \in E$ назовем номер первого нарушенного ограничения или же число $m+1$, если все ограничения в точке выполнены.

Обозначим через $M=M(E)$ максимальное значение индекса в области E и введем вспомогательную задачу

$$g_M^* = g_M(y^*) = \min\{g_M(y) : y \in Y_{M-1}\}, \quad (3.54)$$

где используется $g_{m+1}(y) \equiv f(y)$. Ясно, что если исходная задача будет иметь решение (т.е. $M(E)=m+1$), то оно будет совпадать с решением задачи (3.54). Если же окажется $\nu(y^*)=M(E) < m+1$, то это значит, что исходная задача (3.54) решений не имеет. Согласно индексной схеме вводится *новое понятие результата измерения в точке y* . Считается, что результатом измерения в точке y является вычисление значений

$$\nu(y) \text{ и } z(y) = g_{\nu(y)}(y). \quad (3.55)$$

Пусть выполнено k измерений в точках y^s ($s=1, \dots, k$). Введем множества номеров измерений, выполненных к шагу k и имеющих заданный индекс ν :

$$I_\nu^k = \{i : 1 \leq i \leq k, \nu = \nu(y^i)\}, \quad (3.56)$$

а также множества номеров точек, выполненных измерений, индексы которых больше ν :

$$T_\nu^k = I_{\nu+1}^k \cup \dots \cup I_{m+1}^k \quad (3.57)$$

Определим достигнутые к шагу k модифицированные рекордные значения функции измерений $z(y)$ по группам точек с одинаковым значением индекса ν

$$z_\nu^{*k} = \begin{cases} 0, & T_\nu^k \neq \emptyset \text{ или } I_\nu^k = \emptyset \\ \min\{z(y^s) : s \in I_\nu^k\}, & T_\nu^k = \emptyset. \end{cases} \quad (3.58)$$

В индексной схеме метод многоэкстремальной оптимизации на $k+1$ – м шаге выполняет размещение точки очередного измерения, исходя из следующей вспомогательной цели

$$S_k(y) = g_{\nu(y)}(y) - z_{\nu(y)}^{*k} \rightarrow \min, \quad y \in E \quad (3.59)$$

При этом считается, что методу известны значения функции $S_k(y^s)$ для $s = 1, \dots, k$, которые перерасчитываются по ранее сделанным и запомненным результатам вычислений функций исходной задачи в точках выполненных ранее измерений.

На рисунках 3.10, 3.11 представлены виды функции $S_k(y)$ из (3.59) для задачи, изображенной на рис.3.8, в предположении, что проведено $k=3$ и $k=5$ измерений.

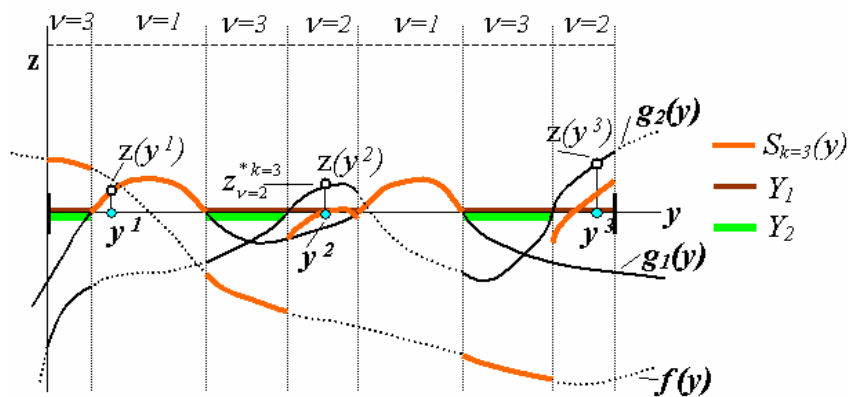


Рис.3.10. Структура вспомогательной функции в индексном методе после $k=3$ измерений

На рис. 3.10 в точках проведенных измерений наибольшее значение индекса ν равно 2, поэтому значения z_{ν}^{*k} для $\nu=1$ и $\nu=3$ равны нулю (т.к. $T_{\nu=1}^k = \emptyset$, $a I_{\nu=3}^k = \emptyset$) и графики функций $g_1(y)$ и $f(y)$ не смещены. Смещение на величину $z_{\nu=2}^{*k}$ происходит только для значений функции $g_2(y)$.

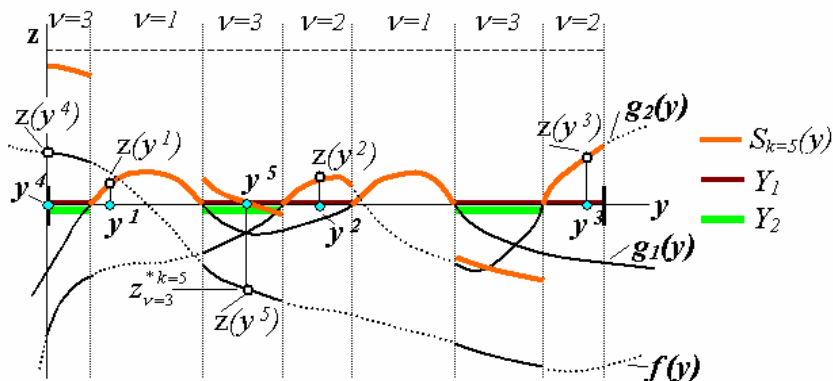


Рис.3.11. Изменение структуры вспомогательной функции после добавления испытания в точке с индексом $\nu=3$

На рис. 3.11 есть два измерения, выполненных в области Y_2 , где определена функция $f(y)$, поэтому наибольшее значение индекса ν в точках измерений равно 3 и, следовательно, значения z_{ν}^{*k} для $\nu=1$ и $\nu=2$ равны нулю, а для $\nu=3$ величина z_{ν}^{*k} отлична от нуля. В результате на этом шаге графики функций $g_1(y)$ и $g_2(y)$ не смещены, а значения функции $f(y)$ смещены на величину $z_{\nu=3}^{*k}$.

Обоснование сходимости индексного метода приведено в [44] применительно к оригинальной авторской версии алгоритма в сочетании с информационно-статистическим методом одномерной многоэкстремальной оптимизации.

Лист регистрации изменений

Дата	Автор	Комментарии
28.06.02	Городецкий С.Ю.	Создание документа
02.07.02	Городецкий С.Ю.	Создание раздела 3.3
04.07.02	Городецкий С.Ю.	Создание раздела 3.4
17.07.02	Городецкий С.Ю.	Внесение изменений в 3.1
18.07.02	Городецкий С.Ю.	Внесение изменений в 3.2
22.07.02	Городецкий С.Ю.	Внесение изменений в 3.3
23.07.02	Городецкий С.Ю.	Внесение изменений в 3.4
12.08.02	Городецкий С.Ю.	Добавлены рисунки в 3.4
25.12.02	Городецкий С.Ю.	Окончательная редакция версии 1
13.09.03	Городецкий С.Ю.	Исправления для версии 2

Глава 4. Математические основы построения и анализа алгоритмов оптимизации

4.1. Модели численных методов оптимизации

4.1.1. Основные обозначения

Предметом интереса настоящей главы является исследование способов решения задач математического программирования, или алгоритмов оптимизации. Прежде, чем переходить к анализу алгоритмов, напомним основные определения, касающиеся класса исследуемых моделей, которые формулируются в виде задач оптимизации.

Пусть $f(y)$ – действительная функция, определенная в области Y N -мерного евклидова пространства R^N и принимающая во всех точках области конечные значения.

Определение 4.1. Задачей оптимизации будем называть задачу следующего вида: найти заданные экстремальные характеристики функции $f(y)$ на множестве Y .

Синонимически данную задачу также часто называют задачей математического программирования.

В зависимости от искомым экстремальных характеристик возможны различные постановки задач оптимизации, обычно связанные с минимальным значением функции

$$f(y^*) = \inf \{f(y) : y \in Y\}, \quad (4.1)$$

и множеством координат глобального минимума

$$Y^* = \text{Arg min} \{f(y) : y \in Y\} = \{y^* \in Y : f(y) = f^*\} \quad (4.2)$$

Символически общую задачу математического программирования будем записывать в виде

$$f(y) \rightarrow \inf, y \in Y. \quad (4.3)$$

и называть также задачей минимизации функции $\varphi(x)$ на множестве Y .

Более детальные формулировки в отношении класса исследуемых задач могут быть найдены в разделе 1.2.

Так как точная верхняя грань функции $f(y)$ на множестве Y

$$\sup \{f(y) : y \in Y\} = -\inf \{-f(y) : y \in Y\} \quad (4.4)$$

то задача определения экстремальных характеристик, связанных с наибольшим значением функции $f(y)$ (задача максимизации), сводится к задаче минимизации функции $-f(y)$. Поэтому везде далее задача математического программирования будет рассматриваться в форме (4.3) и иногда называться просто задачей оптимизации. Функцию $f(y)$ из (4.3) будем называть целевой, минимизируемой или оптимизируемой функцией, множество Y – допустимой областью, а элементы множества Y – допустимыми точками.

Как и ранее, задачу (4.3), для которой заведомо известно, что множество точек глобального минимума Y^* не пусто, будем записывать как

$$f(y) \rightarrow \min, y \in Y. \quad (4.5)$$

Заметим, что часто в литературе задача математического программирования формулируется именно в таком виде.

4.1.2. Формальная модель и общая вычислительная схема

Сформулировав задачу минимизации, мы теперь должны дать ответ на основной вопрос: каким образом ее решать?

Классические подходы математического анализа, использующие возможности аналитического решения, описаны в разделах 2.1 и 2.2. В этом направлении в качестве простого примера рассмотрим следующую задачу. Пусть $f(y)$ - кусочно-гладкая на отрезке $[a,b]$ функция. Тогда минимум $f(y)$ на $[a,b]$ может достигаться лишь в тех точках, где $f'(y) = 0$, либо производная разрывна, либо в граничных точках. Остается найти все такие точки и выбрать из них точку с наименьшим значением. Иными словами, чтобы решить задачу этим способом, требуется:

- а) указание аналитического вида функции;
- б) кусочная гладкость функции;
- в) возможность вычисления производной;
- г) умение решать уравнение $f'(y) = 0$, т.е. задачу поиска корня;
- д) информация о точках разрыва производной или способ определения этих точек.

К сожалению, эти требования на практике выполняются в редчайших случаях. Типичный же случай описывается ситуацией, когда функция $f(y)$ задается алгоритмически, т.е. в виде некоей расчетной схемы, когда по заданному аргументу y рассчитывается значение $f(y)$. В этом случае ни о каких аналитических способах исследования говорить не приходится. Заметим, что даже при аналитическом задании функции и способности посчитать производную исходная задача (4.3) сводится к решению задачи поиска корня, которая по сложности сравнима с задачей минимизации.

Узкая сфера применения аналитических методов обусловила развитие и широкое распространение *численных методов* решения задач оптимизации. Различные формулировки определений численного метода оптимизации даны многими авторами. Общим во всех формулировках является представление метода как некоторой итерационной процедуры, которая (в общем случае последовательно) осуществляет вычисление в точках области поиска определенных характеристик минимизируемой функции (такими характеристиками могут быть значение функции, ее градиента, матрицы вторых производных и т.п.). Назовем операцию вычисления характеристик функции в точке *поисковым испытанием*, а совокупность значений характеристик в этой точке – *результатом испытания*. Далее в настоящей главе в качестве результата испытания будем рассматривать только значение функции в испытываемой точке.

Основываясь на методологии теории исследования операций, дадим формальное определение численного метода оптимизации или, более широко, *модели вычислений* при решении задачи (4.3) [45].

Построение модели вычислений предполагает наличие некоторой априорной (доопытной, имеющейся до начала вычислений) информации о решаемой задаче. Данная информация может быть получена исходя из физической сущности задачи, описывающей моделируемый реальный объект. Такими свойствами могут

быть непрерывность, гладкость, монотонность, выпуклость и т.п. Имеющаяся информация служит для исследователя основанием для отнесения задачи (в нашем случае функции $f(y)$) к тому или иному множеству (классу) Φ . После того, как класс Φ зафиксирован, априорная информация о задаче, используемая исследователем, состоит в том, что ему известна принадлежность задачи к классу Φ .

Следующим важным этапом построения модели вычислений является выбор алгоритма (метода) решения задачи. В самом общем виде численный метод s решения задачи из класса Φ представляет собой набор (кортеж) [7]

$$s = \langle \{G_k\}, \{E_k\}, \{H_k\} \rangle \quad (4.6)$$

в котором семейство функционалов

$$\{G_k\}, \quad k=1,2,\dots \quad (4.7)$$

описывает совокупность правил выбора точек испытаний, последовательность отображений

$$\{E_k\}, \quad k=1,2,\dots \quad (4.8)$$

задает совокупность правил построения приближенного решения (оценки экстремума), а набор

$$\{H_k\}, \quad k=1,2,\dots \quad (4.9)$$

определяет совокупность правил остановки вычислительного процесса.

Порядок проведения испытаний, или *вычислительная схема* алгоритма состоит в следующем.

1. Выбирается точка первого испытания

$$y^1 = G_1(\Phi) \in D \quad (4.10)$$

2. Пусть выбрана точка k -го испытания $y^k \in D$ ($k \geq 1$). Производится вычисление значения функции $z^k = f(y^k)$. После этого имеем следующую поисковую (апостериорную) информацию о функции f :

$$\omega_k = \{(y^1, z^1), (y^2, z^2), \dots, (y^k, z^k)\} \quad (4.11)$$

Полученная информация позволяет сузить класс, которому принадлежит функция $f(y)$ до множества

$$\Phi(\omega_k) = \{\psi \in \Phi : \psi(y^i) = z^i, 1 \leq i \leq k\} \quad (4.12)$$

3. Определяется текущая оценка экстремума (приближенное решение)

$$e^k = E_k(\Phi, \omega_k) \quad (4.13)$$

4. Вычисляется точка очередного испытания

$$y^{k+1} = G_{k+1}(\Phi, \omega_k) \quad (4.14)$$

5. Определяется величина

$$h^k = H_k(\Phi, \omega_k) \in \{0, 1\} \quad (4.15)$$

принимающая одно из двух возможных значений: ноль или единица. Если $h^k = 1$, номер шага поиска k увеличиваем на единицу и переходим к выполнению пункта 2 схемы. Если $h^k = 0$, вычисления прекращаем и в качестве решения задачи берем оценку e^k .

Общая модель вычислений описана.

Пример. Рассмотрим простейший метод решения задачи (4.5) на отрезке $[a, b]$ – метод перебора значений по узлам равномерной сетки. Метод состоит в том, что отрезок разбивается на n равных частей, в точках (узлах) разбиения, в число которых входят и концы отрезка, вычисляются значения функции и в качестве решения задачи рассматривается наименьшее вычисленное значение (и его координата, если это требуется соответствующей постановкой). Для применимости метода достаточно вычислимости функции в любой точке области поиска, так что в качестве априорного класса Φ можно рассмотреть класс функций, определенных на отрезке $[a, b]$ и вычисляемых в каждой его точке.

Для данного метода

$$G_1(\Phi) = a, \quad G_{k+1}(\Phi, \omega_k) = a + k \frac{b-a}{n}, \quad k \geq 1 \quad (4.16)$$

$$H_k(\Phi, \omega_k) = \begin{cases} 0, & k = n+1 \\ 1, & k < n+1 \end{cases} \quad (4.17)$$

$$e^k = f_k^*, \quad (4.18)$$

где

$$f_k^* = \min_{1 \leq i \leq k} f(y^i), \quad (4.19)$$

либо

$$e^k = (f_k^*, y_k^*), \quad (4.20)$$

где

$$y_k^* = \arg \min_{1 \leq i \leq k} f(y^i). \quad (4.21)$$

Контрольные вопросы и упражнения:

1. Предположим, что класс Φ представляет собой класс линейных функций одной переменной и пусть проведено единственное испытание в точке y^1 , обеспечившее результат z^1 . Что из себя будет представлять класс $\Phi(\omega_1)$?
2. Если затем проведено второе испытание $y^2 > y^1$ с результатом z^2 , как выглядит класс $\Phi(\omega_2)$?
3. Пусть для той же задачи после третьего испытания в точке $y^1 < y^3 < y^2$ получено значение z^3 такое, что $z^3 < z^1$, $z^3 < z^2$, какие выводы можно сделать?

4.1.3. Сходимость и оценки решения

Итак, решая задачу минимизации функции $f(y)$, метод поиска порождает (сопоставляет функции) последовательность $y^1, y^2, \dots, y^k, \dots$ координат испытаний, или просто последовательность испытаний (y^k – координата k -го испытания), а также последовательность $z^1, z^2, \dots, z^k, \dots$ результатов испытаний (напомним, что мы ограничились случаем значения функции в качестве результата, т.е. $z^i = f(y^i)$). При этом свойства метода определяются свойствами

последовательностей $\{y^k\}$ и $\{z^k\}$, поэтому исследование метода поиска может быть проведено посредством изучения последовательностей испытаний, им порождаемых.

В связи с этим зададимся вопросом: какие требования должны быть предъявлены к последовательности испытаний численного метода оптимизации? Разумеется, основное требование заключается в том, что проведение испытаний в точках $\{y^k\}$ должно обеспечить на основе результатов $\{z^k\}$ решение задачи, т.е. отыскание решения, соответствующего выбранной постановке. При этом, поскольку вычислитель может осуществить лишь конечное число испытаний, желательно получить точное решение, построив конечную последовательность $\{y^k\}$. К сожалению, такая приятная ситуация имеет место лишь в редких и достаточно простых случаях, например, в задачах линейного программирования. Поэтому часто интересуются асимптотически точной оценкой, рассматривая бесконечную последовательность испытаний (в модели (4.9) $H_k = 1$ для любого k , т.е. условие останова отсутствует) и требуя, чтобы эта последовательность сходилась к точному решению задачи. Поскольку в постановках С-D искомое решение может содержать *несколько* точек минимума, сходимость метода будем понимать в смысле следующего определения.

Определение 4.2. Последовательность испытаний $\{y^k\}$ сходится к решению задачи (4.3), определенному соответствующей постановкой, если:

1) она содержит подпоследовательность $\{\bar{y}^k\}$, для которой

$$\lim_{k \rightarrow \infty} f(\bar{y}^k) = f^* ;$$

2) в случае, когда решение включает одну или несколько точек минимума, для каждой такой точки существует сходящаяся к ней подпоследовательность последовательности $\{y^k\}$.

Последовательность испытаний, сходящуюся к точному решению задачи, будем называть *минимизирующей* последовательностью (данный термин впервые введен в [3]).

Вопросам сходимости в теории методов поиска экстремума уделяется значительное внимание, поскольку асимптотика обеспечивает потенциальную возможность получения точного решения с любой наперед заданной точностью за конечное число испытаний. Но самой по себе такой возможности для практической реализации методов недостаточно. Необходимо еще уметь определять меру близости получаемого приближенного решения к точному решению, т.е. уметь оценивать погрешность решения задачи при конечном числе испытаний.

Рассмотрим следующий простой пример. Пусть класс Φ - класс непрерывных одномерных функций, заданных на отрезке, т.е. априорно известно, что минимизируемая функция $f(y)$ непрерывна в области $Y=[a,b]$. Предположим, что вычислены значения функции $f(y)$ в конечном числе точек y^1, y^2, \dots, y^k . Что после этого можно сказать о координате глобального минимума? Каковы бы ни были точки y^1, y^2, \dots, y^k и значения z^1, z^2, \dots, z^k , для любой точки $y^* \in Y$ ($y^* \notin \{y^1, \dots, y^k\}$) всегда можно построить непрерывную функцию, проходящую

через точки $(y^i, z^i), 1 \leq i \leq k$, т.е. принадлежащую классу $\Phi(\omega_k)$ из (4.12), которая имеет глобальный минимум в точке y^* с любым наперед заданным значением $f^* < \min_{1 \leq i \leq k} z^i$.

Например, в качестве такой функции можно взять интерполяционный полином k -й степени, проходящий через точки $(y^i, z^i), 1 \leq i \leq k$, и точку (y^*, f^*) .

Все сказанное означает, что по результатам конечного числа испытаний никаких выводов о расположении координаты глобального минимума сделать нельзя. Точно так же о величине f^* глобального минимума можно лишь сказать, что

$$f^* \leq f_k^*, \quad (4.22)$$

где f_k^* из (4.19), однако оценить величину

$$\varepsilon_k = |f^* - f_k^*|, \quad (4.23)$$

т.е. погрешность решения задачи, невозможно.

Возможность получения оценок экстремума по конечному числу испытаний зависит от свойств класса функций, которому принадлежит минимизируемая функция, или, другими словами, от априорной информации о функции $f(y)$.

Для примера рассмотрим класс $\Phi = U[a, b]$ одномерных *строго унимодальных* на отрезке $[a, b]$ функций (см. раздел 1.4), т.е. функций, для каждой из которых существует точка $y^* \in [a, b]$ такая, что на отрезке $[a, y^*]$ функция строго убывает, а на отрезке $[y^*, b]$ - строго возрастает. Пусть проведены испытания в точках y^1 и y^2 интервала (a, b) и получены значения целевой функции z^1 и z^2 . Предположим, что $y^1 < y^2$ и $z^1 < z^2$. Тогда в силу унимодальности очевидно, что на отрезке $[y^2, b]$ точка минимума x^* находится не может, и в качестве области локализации координаты минимума можно рассмотреть полуинтервал $[a, y^2)$, называемый *интервалом неопределенности*.

Пусть теперь в общем случае проведено k испытаний в точках $y^1, y^2, \dots, y^k \in (a, b)$ и получены значения z^1, z^2, \dots, z^k . Перенумеруем точки испытаний нижним индексом в порядке возрастания координаты, добавив к ним также концы отрезка поиска a и b , т.е.

$$a = y_0 < y_1 < y_2 < \dots < y_k < y_{k+1} = b \quad (4.24)$$

Тогда интервалом неопределенности будет интервал (y_{i-1}, y_{i+1}) , где номер i определяется из условия $y_i = y_k^*$, где y_k^* из (4.21) (в случаях $i=1$ и $i=k$ интервалами неопределенности будут полуинтервалы $[a, y_2)$ и $(y_{k-1}, b]$ соответственно). Иными словами, для строго унимодальной функции можно построить оценку координаты глобального минимума в виде интервала неопределенности и тем самым оценить погрешность решения задачи (по координате) величиной $\max\{y_{i+1} - y_i, y_i - y_{i-1}\}$, ибо

$$|y_k^* - y^*| < \varepsilon = \max\{y_{i+1} - y_i, y_i - y_{i-1}\} \quad (4.25)$$

Что касается величины глобального минимума, то строгой унимодальности для получения оценки (4.23) недостаточно и требуются более жесткие условия для ее реализуемости.

Другим важным классом функций, допускающим построение оценок экстремума по конечному числу испытаний, является класс функций, удовлетворяющих условию Липшица (см. п.1.4.2.1). Некоторые схемы такого оценивания будут приведены далее.

4.2. Принципы построения методов оптимизации

После того, как решен вопрос о принципиальной возможности построения оценки искомого решения, возникает естественный интерес к исследованию эффективности алгоритма. Хотя понятие эффективности может формулироваться по-разному (например, в терминах скорости сходимости, плотности испытаний и т.п.), тем не менее в любом случае это понятие связывает [7] затраты на поиск с некоторой мерой близости оценки e^k из (4.13) к точному решению задачи.

Приведем формальную постановку задачи определения эффективности алгоритма оптимизации [7, 11, 45]. Согласно этой постановке рассматривается некоторый класс S алгоритмов $s \in S$, предназначенных для решения задач минимизации (4.3) функций f из класса Φ . Вводится вещественная функция $V(f,s)$, называемая критерием эффективности, которая количественно характеризует эффективность решения задачи минимизации функции $f \in \Phi$ с помощью метода $s \in S$. Для определенности будем считать, что чем меньше величина $V(f,s)$, тем выше эффективность алгоритма.

Следуя общей схеме теории исследования операций, в качестве эффективности метода s на классе примем либо гарантированный результат

$$W(s) = \sup_{f \in \Phi} V(f, s) \quad (4.26)$$

(наихудшее возможное "достижение" алгоритма на классе) Φ , либо среднюю по классу Φ эффективность

$$W(s) = \int_{\Phi} V(f, s) dP(f), \quad (4.27)$$

где $P(f)$ — некоторое распределение вероятностей, заданных на классе измеримых подмножеств множества Φ .

Определив для каждого алгоритма $s \in S$ его эффективность $W(s)$, можно поставить задачу нахождения оптимального по данному критерию W метода $s^* \in S$. *Оптимальный алгоритм* s^* должен удовлетворять условию

$$W(s^*) = \inf_{s \in S} W(s) \quad (4.28)$$

Если точная нижняя грань в (4.28) не реализуема (s^* не существует), можно рассмотреть ε -оптимальный алгоритм s^*_ε , для которого (при фиксированном $\varepsilon > 0$) имеет место

$$W(s^*_\varepsilon) \leq \inf_{s \in S} W(s) + \varepsilon \quad (4.29)$$

Если следовать терминологии исследования операций, то методы оптимизации называют также стратегиями оптимизации, а оптимальные (ε -оптимальные) алгоритмы в ситуации гарантированного результата (4.26) –

оптимальными (ε -оптимальными) минимаксными стратегиями. При рассмотрении оптимальности "в среднем", а именно критерия (4.27), который в теории статистических решений называют функцией риска, стратегия, минимизирующая риск, называется байесовской, что приводит к использованию термина "*байесовские методы оптимизации*".

Проблема отыскания оптимального алгоритма решения экстремальной задачи (4.3) в свою очередь сводится к решению экстремальной задачи (4.28), причем, как правило, существенно более сложной, поскольку областями поиска являются нечисловые множества. Как указано в [7], решение этой задачи возможно лишь при наличии соответствующего математического аппарата исследования функции $V(f,s)$, которая, будучи мерой эффективности решения задачи (4.3), неразрывно связана с построением оценок экстремума.

Ранее уже отмечалось, что построение таких оценок определяется свойствами функций, задаваемыми описанием класса Φ . К настоящему времени практически известны лишь два широких класса функций, для которых существует развитый аппарат получения оценок экстремума и оптимальных в том или ином смысле алгоритмов. Одним из этих классов является класс $U[a,b]$ унимодальных функций одной переменной. Именно для этого класса Дж.Кифером впервые была поставлена задача отыскания оптимального алгоритма поиска экстремума и построен минимаксный ε -оптимальный алгоритм – метод Фибоначчи. Дальнейшим исследованиям в этом направлении посвящено много работ, обзоры которых можно найти в [3, 45].

Приведем примеры использования введенных понятий оптимальности применительно к классу $U[a,b]$.

Рассмотрим класс S_K алгоритмов, которые осуществляют ровно K испытаний минимизируемой функции. В качестве погрешности решения задачи минимизации функции $f \in U$ методом $s \in S_K$ введем величину (4.25), т.е. примем

$$V(f,s) = \max\{y_{i+1} - y_i, y_i - y_{i-1}\},$$

где координаты испытаний и концы отрезка перенумерованы нижним индексом в соответствии с (4.24), а номер i определяется из условия $f(y_i) = \min_{1 \leq j \leq K} f(y_j)$.

Выделим в классе S_K подкласс P_K *пассивных* методов поиска. Алгоритм $p \in S_K$ называется пассивным, если решающие правила G_k , $1 \leq k \leq K$, не зависят от результатов проведенных испытаний (значений целевой функции), т.е. $G_k = G_k(y^1, \dots, y^{k-1})$. Заметим, что в этом случае порядок проведения испытаний y^1, \dots, y^K безразличен, поэтому все методы, у которых одинаковый набор точек y^1, \dots, y^K , будем отождествлять.

Поскольку точка глобального минимума y^* неизвестна, то после применения определенного метода $p \in P_K$ любой из интервалов (y_{i-1}, y_{i+1}) , $1 \leq i \leq K$, может оказаться интервалом неопределенности, откуда

$$W(p) = \max_{\varphi \in U} V(f, p) = \max_{1 \leq i \leq K} \max\{y_{i+1} - y_i, y_i - y_{i-1}\} = \max_{0 \leq i \leq K} y_{i+1} - y_i.$$

Минимизировать критерий $W(p)$ в классе P_K будет метод p^* , для которого

$$y^k = a + k \frac{b-a}{K+1}, \quad 1 \leq k \leq K,$$

т.е. метод перебора по равномерной сетке. Действительно, для этого метода $W(p^*) = \frac{b-a}{K+1}$, а для любого отличного от него пассивного алгоритма p всегда найдется интервал (y_{i-1}, y_i) длины $y_i - y_{i-1} > \frac{b-a}{K+1}$. Поэтому $W(p) > W(p^*)$, т.е. метод p является оптимальным.

Если рассмотреть теперь весь класс S_K , т.е. принять во внимание также алгоритмы, учитывающие при планировании поиска полученные значения функции, то для введенной погрешности $V(f, s)$ существует единственный оптимальный алгоритм, именуемый *методом Фибоначчи* F_K . Данный метод основан на использовании числовой последовательности Фибоначчи $u_1, u_2, \dots, u_k, \dots$, для которой

$$u_1 = u_2 = 1, u_k = u_{k-2} + u_{k-1}, k \geq 3.$$

Метод Фибоначчи обеспечивает погрешность

$$W(F_K) = \frac{b-a}{u_{K+2}},$$

что существенно лучше погрешности оптимального пассивного алгоритма $\frac{b-a}{K+1}$ (метод Фибоначчи эффективнее метода перебора в $\frac{u_{K+2}}{K+1}$ раз), причем это превосходство увеличивается с ростом заданного количества испытаний K .

Другим классом, допускающим вывод оптимальных алгоритмов, является класс функций, удовлетворяющих условию Липшица в некоторой метрике. Функции этого класса в общем случае многоэкстремальны. При минимаксном подходе к проблеме конструирования оптимальных методов оптимизации установлена связь этой проблемы с задачей построения оптимального покрытия области поиска [23, 45]. Оптимальные минимаксные алгоритмы построены в работах [12, 45].

Класс $\Phi = Lip[a, b]$ одномерных функций, удовлетворяющих на отрезке $[a, b]$ условию Липшица с константой $L > 0$, рассмотрим в качестве второго примера построения оптимальных алгоритмов.

Пусть в соответствии с решающим правилом алгоритма $s \in S_K$ осуществлено k испытаний и получена информация ω_k . Если упорядочить точки испытаний в соответствии с (4.24), то для каждой точки $y_i, 1 \leq i \leq k$, можно построить функцию

$$\psi_i(y) = f(y_i) - L|y - y_i|, y \in [a, b],$$

для которой вследствие условия Липшица

$$f(y) \geq \psi_i(y), y \in [a, b],$$

откуда

$$f(y) \geq f_k^-(y) = \max \{ \psi_i(y), 1 \leq i \leq k \}, y \in [a, b]$$

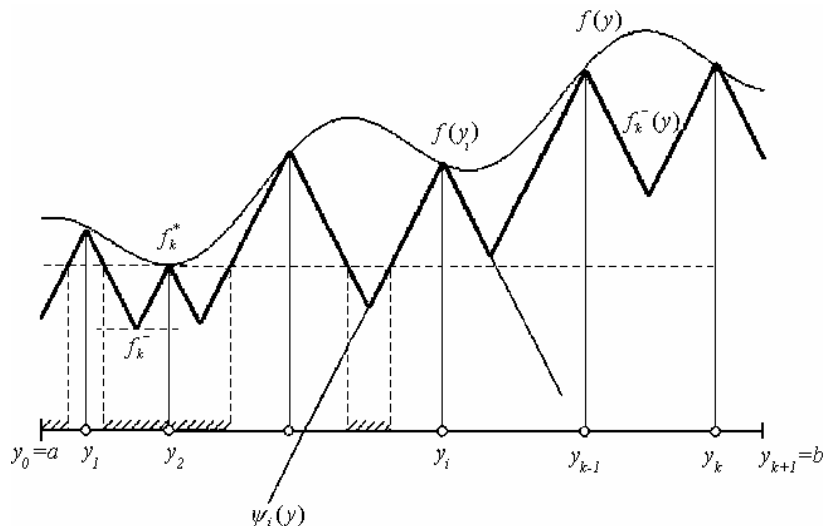


Рис.4.1. Миноранта и оценки решения по результатам измерения липшицевой функции

Функция $f_k^-(y)$ является кусочно-линейной *минорантой* для целевой функции $f(y)$, построенной по результатам испытаний ω_k . Нетрудно видеть, что глобально-минимальное значение f^* целевой функции находится в границах

$$f_k^- \leq f^* \leq f_k^*,$$

где

$$f_k^- = \min \{f_k^-(y) : y \in [a, b]\},$$

$$f_k^* = \min \{z_i = f(y_i) : 1 \leq i \leq k\}.$$

Аналогично можно утверждать, что любая точка глобального минимума y^* принадлежит множеству (на рисунке отмечено штриховкой)

$$D_k = \{y \in [a, b] : f_k^-(y) \leq f_k^*\}.$$

Таким образом, для липшицевых функций возможна как оценка глобально-минимального значения f^* , если положить $e^k = [f_k^-, f_k^*]$, так и оценка координаты y^* глобального минимума, если принять $e^k = D_k$. Рассмотрим эти случаи.

Вначале допустим, что принята оценка по координате. Положим эффективность $V(f, s)$ применения алгоритма $s \in S_K$ для минимизации функции $f \in Lip[a, b]$ равной мере множества D_K , которая равна суммарной длине интервалов, его составляющих. В этом случае любому методу $s \in S_K$ соответствует $W(s) = b - a$, если функция $f(y) = const$. Следовательно, каждый алгоритм $s \in S_K$ удовлетворяет условиям (4.26), (4.28) и является минимаксным. Подобная ситуация, разумеется, не представляет никакого интереса, поэтому выберем другой вариант оценки.

Рассмотрим в качестве критерия эффективности $V(f, s)$ алгоритма $s \in S_K$ при минимизации функции $f \in Lip[a, b]$ величину погрешности оценки минимального значения, т.е. $V(f, s) = f_k^* - f_k^-$.

Перепишем величину f_k^- как

$$f_k^- = \min \{ \psi(y_i^*) : 1 \leq i \leq k+1 \},$$

где

$$\psi(y_i^*) = \min \{ f_k^-(y) : y \in [y_{i-1}, y_i] \}, 1 \leq i \leq k+1.$$

Из условия Липшица и определения миноранты $f_k^-(y)$

$$y_i^* = (y_{i-1} + y_i) / 2 - (z_i - z_{i-1}) / (2L), 2 \leq i \leq k,$$

$$y_1^* = a, y_{k+1}^* = b$$

откуда

$$\psi(y_i^*) = \begin{cases} (z_i + z_{i-1} - L\Delta_i) / 2, & 2 \leq i \leq k \\ z_i - L\Delta_i, & i = 1 \\ z_{i-1} - L\Delta_i, & i = k+1 \end{cases}$$

где $\Delta_i = y_i - y_{i-1}, 1 \leq i \leq k+1$.

Погрешность определения минимального значения целевой функции на отрезке $[y_{i-1}, y_i]$ не превосходит величины

$$\delta_i = \begin{cases} \min \{ z_{i-1}, z_i \} - \psi(y_i^*) \leq L\Delta_i / 2, & 2 \leq i \leq k \\ z_1 - \psi(y_1^*) \leq L\Delta_1, & i = 1 \\ z_k - \psi(y_{k+1}^*) \leq L\Delta_{k+1}, & i = k+1 \end{cases}$$

Тогда

$$f_k^* - f_k^- \leq \delta = \max \{ \delta_i : 1 \leq i \leq k+1 \}.$$

Рассмотрим в классе S_K метод s^* перебора по узлам сетки

$$y^i = y_i = a + (2i-1) \frac{(b-a)}{2K}, 1 \leq i \leq K.$$

Для данного метода $\Delta_1 = \Delta_{k+1} = \frac{b-a}{2K}$, $\Delta_i = \frac{b-a}{K}, 1 \leq i \leq K$, поэтому для любой функции $f(y)$

$$V(f, s^*) = f_k^* - f_k^- \leq \frac{L(b-a)}{2K}.$$

С другой стороны, если $f(y) = const$, то

$$\delta_i = \begin{cases} L\Delta_i / 2, & 2 \leq i \leq k, \\ L\Delta_i, & i = 1, \\ L\Delta_i, & i = k+1. \end{cases}$$

и, следовательно, $V(f, s^*) = \frac{L(b-a)}{2K}$. Но тогда согласно (4.26) функция-константа

является наихудшей и $W(s^*) = \frac{L(b-a)}{2K}$.

Для произвольного метода s либо $y_1 - a \geq \frac{b-a}{2K}$, либо $b - y_K \geq \frac{b-a}{2K}$, либо найдется такой интервал (y_{i-1}, y_i) , $2 \leq i \leq K$, что $y_i - y_{i-1} \geq \frac{b-a}{K}$, откуда для функции $f(y) = \text{const}$ имеем $V(f, s) \geq \frac{L(b-a)}{2K}$, т.е. $W(s) \geq W(s^*)$, что доказывает оптимальность метода перебора.

Что касается байесовых методов поиска экстремума, то Й.Б.Мощкусом [32] показано, что задача их построения может быть сведена к решению некоторой системы рекуррентных уравнений. Однако и в этом случае задача остается настолько сложной, что оптимальные байесовы методы точно не реализуемы.

Формулировка критериев оптимальности в виде (4.26), (4.27) ориентирована на использование только априорной информации и не учитывает одной важной особенности организации вычислительного процесса поиска минимума, а именно, возможности учета получаемой в процессе поиска (апостериорной) информации о функции. Это приводит для критерия (4.26) к ориентации на худший, а для критерия (4.27) – на "типичный" (наиболее вероятный) случай, хотя оптимизируемая функция может существенно отличаться от худшей, либо типичной. При минимаксном подходе в классе липшицевых функций оказывается, что вследствие ориентации на худший случай, которым является функция-константа, многие оптимальные алгоритмы представляют собой метод перебора по равномерной сетке. В связи с этим А.Г.Сухарев ввел понятие *последовательно-оптимального* алгоритма как метода, который являясь минимаксным, вместе с тем на каждом шаге поиска наилучшим образом использует апостериорную информацию о минимизируемой функции (строгое определение можно найти в [45]). Понятие последовательной оптимальности можно применить и к байесовским методам поиска [7]. Более экономные последовательно-оптимальные алгоритмы удалось построить А.Г.Сухареву для некоторых подклассов липшицевых функций.

4.3. Одношагово-оптимальные методы оптимизации

Трудности конструирования оптимальных в смысле (4.28) алгоритмов поиска привели к использованию более простых понятий оптимальности. Одним из таких понятий является так называемая *одношаговая оптимальность* [7, 9, 11, 45], когда алгоритм размещает очередное испытание наилучшим образом, предполагая, что оно является последним.

4.3.1. Принцип одношаговой оптимальности

Введение понятия одношаговой оптимальности связано не только с необходимостью упрощения формулировки принципа оптимальности. Довольно часто в процессе поиска приходится уточнять модель решаемой задачи (такая ситуация имеет место для многих вычислительных проблем). Поэтому на каждом шаге поиска k справедливо свое предположение $f \in \Phi_k$, меняющееся от шага к шагу. В этом случае использование принципа одношаговой оптимальности совершенно естественно.

Приведем формальное описание *одношагово-оптимального метода поиска экстремума* [7, 9]. С этой целью вернемся к формальной схеме метода оптимизации (4.9) и для любого возможного ω_k из (4.11) определим класс $\Phi(\omega_k)$

из (4.12), т.е. множество функций из Φ , допускающих реализацию данного ω_k . В исходном множестве алгоритмов S выделим подкласс $S(\omega_k)$ таких методов, которые для любой функции $f \in \Phi(\omega_k)$ после первых k испытаний реализуют данное ω_k . Введем последовательность функций $V_{k+1}(\omega_k, y, z)$, определяющих эффективность проведения очередного испытания в точке y и получения результата $z=f(y)$ для $f \in \Phi(\omega_k)$.

Тогда алгоритм s может быть описан набором

$$s = \langle \{G_k\}, \{E_k\}, \{H_k\}, \{V_k\} \rangle \quad (4.30)$$

Множество таких алгоритмов обозначим через S_0 . Введем величину

$$W_{k+1}(y^{k+1}) = \sup_{f \in \Phi(\omega_k)} V_{k+1}(\omega_k, y^{k+1}, f(y^{k+1})) \quad (4.31)$$

называемую функцией гарантированного результата. Выбор операции sup/inf определяется тем, как введена функция эффективности (случай sup соответствует варианту "чем меньше $V_{k+1}(\omega_k, y, z)$, тем лучше").

Тогда метод $s \in S_0$ называется (минимаксным/максиминным) одношагово-оптимальным, если точки испытаний y^k , $k=1,2,\dots$, порождаемые методом s , удовлетворяют условиям

$$W_k(y^k) = \min_{\tilde{y}^k \in Y} W_k(\tilde{y}^k) \quad (4.32)$$

Если вместо гарантированного результата (4.31) рассмотреть математическое ожидание

$$W_{k+1}(y^{k+1}) = \int_{\Phi(\omega_k)} V_{k+1}(\varphi, y^{k+1}, z^{k+1}) dP(z^{k+1} / \omega_k, y^{k+1}), \quad (4.33)$$

определяющее функцию средней эффективности, где $P(z^{k+1} / \omega_k, y^{k+1})$ - условное по отношению к результатам испытаний распределение вероятностей, то метод, определяемый соотношением (4.32), будет называться байесовским одношагово-оптимальным алгоритмом, или методом, одношагово-оптимальным в среднем.

Приведем примеры одношагово-оптимальных алгоритмов.

4.3.2. Метод ломаных как одношагово-оптимальный алгоритм

Классическим примером минимаксного одношагово-оптимального алгоритма является метод С.А.Пиявского. Данный метод предполагает, что минимизируемая функция одного аргумента $f(y)$ удовлетворяет на отрезке $[a, b]$ условию Липшица с константой $L > 0$, т.е. $f(y) \in Lip[a, b]$.

Идея метода геометрически проста: предлагается проводить очередное испытание в точке минимума кусочно-линейной миноранты $f_k^-(y)$. Можно показать, что такой выбор точки очередного испытания реализует принцип одношаговой оптимальности, если в качестве функции эффективности за один шаг выбрать функцию

$$V_{k+1}(\omega_k, y, z) = \min \{ f_k^*(z); z \} - \min_{u \in [a, b]} \{ f_k^-(u); z - L|y - u| \},$$

которая описывает взятое с обратным знаком ожидаемое уменьшение погрешности на $k+1$ -м шаге поиска по сравнению с погрешностью на предыдущем k -м шаге, т.е. выбор координаты очередного испытания в точке минимума миноранты обеспечивает максимальное уменьшение текущей погрешности.

4.3.3. Информационно-статистический метод Р.Г.Стронгина

В качестве первого примера построения одношагово-оптимальных байесовских алгоритмов рассмотрим информационно-статистический алгоритм глобального поиска, предложенный Р.Г.Стронгиным. В основе алгоритма лежит вероятностная модель функций, заданных на конечном множестве точек.

Итак, пусть на отрезке $[a, b]$ определена функция $f(y)$. Будем минимизировать данную функцию на конечной равномерной сети точек

$$a = y(0) < y(1) < y(2) < \dots < y(n-1) < b = y(n)$$

избранного отрезка. Поскольку любая функция $f(y(i)), 0 \leq i \leq n$, полностью определяется значениями $f_i = f(y(i)), 0 \leq i \leq n$, то ее можно представить как точку

$$f = (f_0, \dots, f_n) \in R^{n+1}$$

$(n+1)$ -мерного евклидова пространства R^{n+1} . Тогда возможный способ описания априорных предположений о минимизируемой функции $f \in R^{n+1}$ состоит в задании плотности $\varphi(f)$ распределения вероятностей на пространстве R^{n+1} .

Если плотность $\varphi(f)$ положительна и непрерывна, то с вероятностью единица минимум оптимизируемой функции достигается в единственной точке, поэтому для его оценки достаточно подсчитать вероятности $\eta(\alpha)$ расположения минимума в точках α , совпадающих с узлами сетки. Более того, после проведения испытаний в некоторых узлах сетки и получения поисковой информации ω_k возможно получение условных (апостериорных) вероятностей $\eta(\alpha / \omega_k)$, обеспечивающих текущие оценки экстремума.

К сожалению, вычисление данных вероятностей требует интегрирования по сложным многомерным областям, поэтому предлагается аппарат приближенного описания апостериорных вероятностей $\eta(\alpha / \omega_k)$ с помощью введения дополнительной дискретной случайной величины α , распределение вероятностей $\xi(\alpha)$ которой на узлах сетки связано с априорным вероятностным описанием $\varphi(f)$ с помощью разложения

$$\varphi(f) = \sum_{\alpha \in Y_n} \varphi(f / \alpha) \xi(\alpha)$$

где Y_n - множество узлов сетки. Введение такого разложения позволяет установить основанные на методе малого параметра условия близости апостериорных вероятностей $\eta(\alpha / \omega_k)$ и $\xi(\alpha / \omega_k)$ и тем самым перейти к вычислению условных вероятностей состояния природы, которые существенно проще.

В качестве конкретного примера применения информационно-статистического подхода рассмотрим вероятностное описание, в котором априорная плотность $\varphi(f)$ представлена разложением $\varphi(f / \alpha), 0 \leq \alpha \leq n$, где

$$\varphi(f/\alpha) = \prod_{i=0}^n \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left\{ - \sum_{i=0}^n \frac{(\Delta f_i - m_i(\alpha))^2}{2\sigma_i^2} \right\}.$$

Здесь

$$\begin{aligned} \Delta f_i &= f_i - f_{i-1}, 1 \leq i \leq n, \quad \Delta f_0 = f_0, \\ \sigma_i &= cm, 1 \leq i \leq n, \quad c > 0, \quad m > 0, \\ m_i(\alpha) &= \begin{cases} -m, & i \leq \alpha, \\ m, & i > \alpha, \end{cases} \quad 1 \leq i \leq n. \end{aligned}$$

Согласно данному описанию первые разности Δf_i функции и значение f_0 этой функции в точке $y(0)$ при любом фиксированном α являются реализациями независимых случайных величин, подчиненных нормальному закону со стандартами $\sigma_i = cm, 1 \leq i \leq n, \sigma_0$ и математическими ожиданиями $m_i(\alpha), 1 \leq i \leq n, m_0$.

Принятые предположения являются некоторым вероятностным аналогом условия Липшица, которое ограничивает первые разности минимизируемой функции и обеспечивает ее равномерную непрерывность. В самом деле, согласно описанию дисперсии разностей Δf_i функции f ограничены, и, кроме того, эти разности являются реализациями независимых нормальных величин, что при $n \rightarrow \infty$ и $mn = const$ обеспечивает вероятностное свойство, аналогичное равномерной непрерывности.

Кроме того, априорное математическое ожидание

$$\mu_\alpha(i) = m_0 - m \times \begin{cases} i, & i \leq \alpha \\ 2\alpha - i, & i > \alpha \end{cases}$$

значений f_i функции f при фиксированном α имеет единственный локальный минимум, т.е. функции, близкие к константе, являются маловероятными.

Пусть теперь проведены испытания в точках сетки

$$a = y_0 < y_1 < y_2 < \dots < y_k < b = y_{k+1}$$

и получена информация ω_k . Тогда апостериорные вероятности состояния природы могут быть выражены как

$$\xi(\alpha/\omega_k) = B \cdot h(\omega_k, \alpha) \xi(\alpha),$$

где B - константа, а для любого α из отрезка $y_{s-1} \leq \alpha \leq y_s, i_{s-1} \in J, i_s \in J, 1 \leq s \leq k$,

$$h(\omega_k, \alpha) = \exp \left\{ - \frac{(\alpha - \alpha_s^*)^2}{2\rho_s^2} \right\} \exp \left\{ \frac{R(s)}{2mc^2} \right\}$$

$$\alpha_s^* = \frac{y_s - y_{s-1}}{2} - \frac{z_s - z_{s-1}}{2m},$$

$$R(s) = m\Delta_s + \frac{(z_s - z_{s-1})^2}{m\Delta_s} - 2(z_s - z_{s-1}),$$

$$\rho_s^2 = \left(\frac{c}{2}\right)^2 \Delta_s, \quad \Delta_s = y_s - y_{s-1}.$$

Пусть распределение $\xi(\alpha)$ является равномерным ($\xi(\alpha) = \frac{1}{n+1}$), т.е. все исходные состояния природы равновероятны, и параметр m удовлетворяет условию

$$m > \max_{1 \leq s \leq k} \left| \frac{z_s - z_{s-1}}{\Delta_s} \right|$$

тогда при достаточно малом c наивероятнейшее состояние природы достигается в точке интервала (y_{t-1}, y_t) , для которого величина $R(t)$, которую назовем *характеристикой интервала*, максимальна, и совпадает с наивероятнейшей точкой расположения минимума.

Более того, точка α_t^* будет одношагово-оптимальной точкой, если в качестве функций потерь $V_{k+1}(\omega_k, y, z)$, описывающих эффективность приближения, полученного в результате проведенных испытаний, рассмотреть либо функцию

$$V_{k+1}(\omega_k, y, z) = \begin{cases} 0, & f^* = \min\{f_k^*; z\} \\ w, & f^* < \min\{f_k^*; z\} \end{cases}$$

где $w > 0$, $f^* = \min_{0 \leq i \leq n} f_i$, $f_k^* = \min_{0 \leq s \leq k} z_s$, либо

$$V_{k+1}(\omega_k, y, z) = \sum_{s=0}^k z_s.$$

В первом случае потери считаются равными нулю, если значение в точке абсолютного минимума уже вычислено, и положительны, если это не так.

Во втором варианте потери будут тем меньше, чем меньше значения функции в точках испытаний.

Построенный метод легко распространяется на функции с непрерывно изменяющимся аргументом. Для этого достаточно устремить число точек равномерной сетки к бесконечности ($n \rightarrow \infty$) так, чтобы в пределе был получен конечный интервал $D=[a, b]$. При этом полученные расчетные формулы для характеристики интервала $R(s)$ и наилучшей точки измерения на нем α_s^* не изменятся. Координату точки измерения на s -м интервале в дальнейшем вместо α_s^* будем обозначать через \tilde{y}_s .

Заметим, что полученный метод можно записать в другой форме, формально используя в качестве функции средней эффективности выражение вида

$$W_{k+1}(y) = m(y_s - y_{s-1}) \left(1 + \frac{(z_s - z_{s-1})^2}{(m(y_s - y_{s-1}))^2} \right) - 2(z_s + z_{s-1}) - 4(2y - (y_s + y_{s-1})) + \frac{(z_s - z_{s-1})^2}{m} / (y_s - y_{s-1}) \quad (4.34)$$

для $y \in [y_{s-1}; y_s]$. Геометрически этой функции соответствует выпуклая вверх парабола, принимающая на концах интервала значения $-4f_{s-1}$, $-4f_s$.

С использованием (4.34) правило выбора точки очередного испытания в информационно-статистическом методе может быть записано в виде

$$W_{k+1}(y^{k+1}) = \max \{W_{k+1}(y) : y \in D = [a; b]\}.$$

4.3.4. Одношагово-оптимальный байесовский метод Х.Кушнера

В качестве второго примера одношагово-оптимального байесовского алгоритма приведем метод, предложенный Х.Кушнером. В качестве вероятностной модели целевой функции рассматривается гауссовский случайный процесс с независимыми безгранично-делимыми приращениями (такой процесс получил название *винеровского*). В рамках модели любое приращение функции $\varphi(x_2) - \varphi(x_1)$ распределено по нормальному закону с нулевым средним и дисперсией вида $\sigma|x_2 - x_1|$, т.е. пропорциональной приращению аргумента. Их теории случайных процессов известно, что реализации такого процесса непрерывны с вероятностью 1 и нигде не дифференцируемы. Таким образом, данная модель характеризует поведение многоэкстремальной непрерывной негладкой функции. Заметим, что выбор нулевого среднего означает, что до проведения измерений наиболее вероятной, или "типичной" функцией является функция-константа, что существенно отличает данную модель от модели информационно-статистического алгоритма. Результаты, к которым приводит такое различие, будут обсуждаться далее.

Введем функцию оптимальности на шаге следующим образом. Будем считать, что y — точка планируемого $(k+1)$ -го испытания, а $z=f(y)$ — пока еще не известное нам значение целевой функции в этой точке. С учетом планируемого измерения $\omega_{k+1} = \omega_k \cup \{(y; z)\}$. Выберем вид функции эффективности таким образом, чтобы $V_{k+1}(f, \omega_{k+1})$ зависела от функции f только через посредство значения $z=f(y)$, тогда получим $V_{k+1}(f, \omega_{k+1}) = V_{k+1}(\omega_k; z)$. Примем

$$V_{k+1}(\omega_k; z) = \begin{cases} 1, & z \leq f_k^* - \delta_k, \\ 0, & z > f_k^* - \delta_k, \end{cases} \quad (4.35)$$

т.е. потери постоянны и равны единице, если значение функции в точке очередного испытания меньше текущего минимального после k -го шага значения f_k^* с "запасом", который задается величиной параметра $\delta_k > 0$, и равны нулю в противном случае.

При вычислении функции средней эффективности $W_{k+1}(y)$ для винеровской модели поведения целевой функции используется функция условного распределения с гауссовой плотностью $p(z/\omega_k, y)$, вид математического ожидания и дисперсии которой приведен в (1.57), (1.58) главы 1 и проиллюстрирован на рис.1.22. С учетом (4.35) получим

$$W_{k+1}(y) = \int_{-\infty}^{+\infty} V_{k+1}(\omega_k, z) p(z/\omega_k, y) dz = P(f(y) \leq f_k^* - \delta_k / \omega_k).$$

Согласно принципу одношаговой оптимальности в среднем, выбор точки очередного испытания определяется правилом

$$W_{k+1}(y^{k+1}) = \max_{a \leq y \leq b} P(f(y) \leq f_k^* - \delta_k / \omega_k),$$

т.е. y^{k+1} находится из условия максимизации вероятности вычисления в точке y^{k+1} значения целевой функции меньшего достигнутого минимального значения f_k^* по крайней мере на $\delta_k > 0$.

Если предположить, что проведены испытания в точках (4.24), то можно показать, что на каждом отрезке $[x_{i-1}, x_i]$, $1 \leq i \leq k$, максимизация функции вероятности $P(f(y) \leq f_k^* - \delta_k / \omega_k)$ эквивалентна максимизации выражения

$$W_k(\omega_k, y) = (f_k^* - \delta_k - M[f(y)/\omega_k]) / D[f(y)/\omega_k], \quad (4.36)$$

где математическое ожидание и дисперсия определяются соотношениями (1.57), (1.58). Это выражение достигает своего максимума в точке

$$\tilde{y}_i = y_{i-1} + \frac{(y_i - y_{i-1})(f_k^* - \delta_k - f(y_i))}{2(f_k^* - \delta_k) - f(y_i) - f(y_{i-1})}$$

со значением

$$R(i) = - \frac{4(f_k^* - \delta_k - f(y_i))(f_k^* - \delta_k - f(y_{i-1}))}{y_i - y_{i-1}}.$$

Теперь в качестве точки очередного испытания достаточно выбрать ту точку \tilde{y}_i , которой соответствует максимальная величина $R(i)$.

Замечания.

1. При программной реализации метода Х.Кушнера параметр δ_k обязательно должен адаптивно изменяться в ходе поиска решения задачи (как правило, убывать), оставаясь отделенным от нуля некоторым малым положительным значением $0 < \delta \ll 1$: $\delta_k \geq \delta > 0$.

2. Настройка параметра δ_k может выполняться с помощью специального алгоритма, аналогичного приведенному ниже в (4.43). Применение алгоритма настройки δ_k требует оценивания параметра σ винеровской модели, например, с помощью соотношения $\sigma = \gamma \tilde{\sigma}_k$, где $\gamma > 1$,

$$\tilde{\sigma}_k = (1/(b-a)) \sum_{i=2}^k (f(y_i) - f(y_{i-1}))^2.$$

4.3.5. Одношагово–оптимальный метод на основе адаптивных вероятностных моделей для задач с ограничениями

Понятие одношаговой оптимальности естественным образом встраивается в общую теорию конструирования алгоритмов на основе адаптивных вероятностных (стохастических) моделей, введенных в пункте 1.4.2.3 [14, 16, 17, 18].

Рассматриваемый ниже пример построения одношагово–оптимального метода интересен в первую очередь тем, что он во-первых, хорошо иллюстрирует технологию распространения принципов одношаговой оптимальности на задачи более сложные, чем рассмотренные в двух предыдущих разделах. Во-вторых, он показывает нетрадиционный способ учета ограничений в многоэкстремальной оптимизации, дополняющий общие подходы, изложенные в главе 3.

Итак, рассмотрим задачу с ограничениями–неравенствами

$$f(y) \rightarrow \min, y \in Y, f : D \rightarrow R^1, \quad (4.37)$$

$$Y = \{y \in D : g(y) \leq 0\}, g = (g_1, \dots, g_m) : D \rightarrow R^m, \quad (4.38)$$

$$D = \{y \in R^N : a \leq y \leq b\}, N \geq 1, \quad (4.39)$$

предполагая, что функции f и g принадлежат некоторому подклассу кусочно-непрерывных функций. Конкретные требования к ним будут оговорены позднее.

Для того, чтобы при построении метода упростить учет набора ограничений, будем использовать функцию обобщенного ограничения (1.66), введенную в пункте 1.4.2.3 главы 1. Она имеет вид

$$G(y) = \max \{g_1(y), \dots, g_m(y)\}.$$

Предположим, что накоплена поисковая информация ω_k , включающая измерения набора функций $Q(y) = (f(y), g_1(y), \dots, g_m(y))$, описывающего поставленную задачу. Таким образом, имеем $\omega_k = \omega_k(Q) = \omega_k(f, g)$, что позволяет выделить отдельно поисковую информацию, накопленную о функции f , т.е. $\omega_k(f)$, и поисковую информацию о значениях функции g , т.е. $\omega_k(g)$.

Будем считать, что для произвольно взятой точки y из D неизвестные значения $z^f = f(y)$ и $z^G = G(y)$ являются реализациями независимых случайных величин $\xi^f(y)$ и $\xi^G(y)$ с функциями распределения $F^f(z^f/\omega_k(f), y)$ и $F^G(z^G/\omega_k(g), y)$, вид которых определяется соотношениями (1.62)–(1.65) и (1.67)–(1.70) из пункта 1.4.2.3 первой главы. Совокупность этих распределений образует адаптивную вероятностную модель решаемой задачи.

Следует обратить внимание на то, что модель поведения функции G , определяемая функцией распределения F^G , строится не по измерениям функции G , а по результатам $\omega_k(g)$ измерений исходного набора функций ограничений g_1, \dots, g_m . Существующие здесь отличия иллюстрируются на рис.1.25.

Перейдем к построению метода. Прежде чем ввести функцию эффективности, определим текущую оценку решения следующим образом

$$f_k^* = \begin{cases} +\infty, & \text{если } \forall i = 1, \dots, k \exists j \in \{1, \dots, m\} : g_j(y^i) > 0 \\ \min \{f^i : i = 1, \dots, k \wedge g(y^i) \leq 0\}, & \text{если } \exists i \in \{1, \dots, k\} \forall j \in \{1, \dots, m\} : g_j(y^i) \leq 0. \end{cases} \quad (4.40)$$

Эта оценка определяет наименьшее среди допустимых точек измерений вычисленное значение целевой функции. Если допустимых точек к текущему шагу не встретилось, оценке приписывается бесконечно большое значение.

Теперь аналогично (4.35) введем функцию эффективности, зависящую от накопленной поисковой информации $\omega_k(f, g)$, а также от результатов z^f и z^G измерений функций f и G в точке планируемого испытания y .

$$V_{k+1}(\omega_k(f, g), z^f, z^G) = \begin{cases} 1, & \text{при } z^f \leq f_k^* - \delta_k, \quad z^G \leq -\varepsilon^G \\ 0, & \text{в остальных случаях} \end{cases} \quad (4.41)$$

При таком способе определения функции эффективности она принимает положительные значения только в тех случаях, когда точка проведенного измерения оказывается допустимой с «запасом» $\varepsilon^G > 0$, а текущую оценку решения f_k^* удастся улучшить, по крайней мере, на величину $\delta_k > 0$. Функция средней эффективности определяется соотношением

$$W_{k+1}(y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} V_{k+1}(\omega_k(f, g), z^f, z^G) dF^f(z^f/\omega_k(f), y) dF^G(z^G/\omega_k(g), y),$$

которое, как нетрудно видеть, может быть записано в виде

$$W_{k+1}(y) = W_k(\omega_k(f, g), y) = F^f(f_k^* - \delta_k / \omega_k(f), y) F^G(-\varepsilon^G / \omega_k(g), y) = \\ = P(f(y) \leq f_k^* - \delta_k / \omega_k(f)) P(G(y) \leq -\varepsilon^G / \omega_k(g)) \quad (4.42)$$

Согласно принципам одношаговой оптимальности очередное испытание должно проводиться в точке y^{k+1} , для которой

$$W_{k+1}(y^{k+1}) = \max\{W_{k+1}(y) : y \in D\}.$$

В этой точке достигается максимума произведение вероятности вычисления значения $f(y)$, улучшающего f_k^* не менее чем на δ_k , на вероятность выполнения в этой точке ограничений в усиленной форме.

После определения точки y^{k+1} вычисляются значения всех функций $f^{k+1}, g_1^{k+1}, \dots, g_m^{k+1}$ и пополняется набор поисковой информации

$$\omega_{k+1}(f, g) = \omega_k(f, g) \cup (y^{k+1}, f^{k+1}, g^{k+1}),$$

что позволяет аналогично выполнить следующую итерацию.

Построенный метод будет проанализирован в последующих разделах 4.5, 4.6 этой главы с точки зрения его сходимости и особенностей размещения точек испытаний.

Заметим, что входящий в функцию эффективности параметр δ_k существенно влияет на поведение построенного метода. Для его эффективной работы параметр δ_k должен постепенно адаптивно уменьшаться до некоторого малого значения $0 < \delta \leq \varepsilon^f$, где ε^f — заданная точность по значению целевой функции. Возможный алгоритм настройки этого параметра состоит в следующем. Зададим пороговую величину β ($0 < \beta < 0.5$) для значения вероятности $W_{k+1}(y^{k+1})$. Если это значение окажется меньше β , то δ_k уменьшается вдвое (это приведет к увеличению $W_{k+1}(y^{k+1})$).


АЛГОРИТМ НАСТРОЙКИ δ_k .

$$\text{ПОКА } (\max\{W_{k+1}(y) : y \in D\} < \beta) \quad (4.43)$$

ДЕЛАЙ { Если $\delta_k \leq \varepsilon^f$ то *останов поиска*,
 ИНАЧЕ $\delta_k := \delta_k / 2$ }

Использование этого алгоритма требует правильного оценивания параметров σ_k^f и σ_k^G модели задачи (смотри (1.64), (1.70)). Это можно сделать с использованием метода максимального правдоподобия. Например, на шаге k определить численно значение σ_k^f из приближенного решения вспомогательной задачи поиска максимума

$$\prod_{i=k-r+1}^k \frac{\partial F^f}{\partial z^f}(f^i / \omega_{k-r}(f), y^i) \rightarrow \max_{\sigma_k^f \geq 0} \quad (4.44)$$

 **Замечание.** Рассмотренный метод следует рассматривать как некоторую теоретическую схему, поскольку его нельзя точно вычислительно реализовать. Однако он допускает приближенные вычислительные реализации с использованием либо приближенных методов максимизации функции (4.42) либо компонентного подхода, примеры применения которого рассматриваются в разделе 5.7.

4.3.6. Асимптотическая оптимальность

Дальнейшим упрощением принципа оптимальности является рассмотрение *асимптотически оптимальных* алгоритмов. Для определения этого понятия обозначим через S_K множество методов (4.9), в которых остановка осуществляется ровно через K шагов поиска. Введем величину

$$W(K) = \inf_{s \in S_K} W(s), \quad (4.45)$$

где $W(s)$ из (4.26) или (4.27) и предполагается, для определенности, что уменьшение $W(s)$ соответствует улучшению качества метода s . Для оптимального метода $s_K^* \in S_K$ очевидно выполняется $W(s_K^*) = W(K)$.

Алгоритм \hat{s}_K называется асимптотически оптимальным, если $W(\hat{s}_K)/W(K) \rightarrow 1$ при $K \rightarrow \infty$.

В заключение отметим, что можно выстроить условную иерархию "сложности" рассмотренных принципов оптимальности. Так, самым сложным и порождающим наиболее эффективные алгоритмы является принцип последовательной оптимальности, полнее всего учитывающий информационную составляющую процесса оптимизации. Вслед за ним по сложности построения можно поставить оптимальность согласно (4.28) (или ε -оптимальность (4.29)), ориентированную только на априорное знание. Далее следует принцип оптимальности на один шаг вперед и, наконец, асимптотическая оптимальность.

4.4. Теоретические основы сходимости одномерных алгоритмов глобального поиска

Ранее мы уже отмечали, что важнейшим свойством численного метода является его сходимость к искомому решению, в нашем случае – к глобальному минимуму целевой функции задачи (4.6). При этом характер сходимости во многом определяет эффективность метода оптимизации. Настоящий параграф посвящен рассмотрению с единых теоретических позиций вопросов сходимости для широкого класса численных методов поиска глобального экстремума функций одной переменной, называемого классом характеристических алгоритмов и включающего многие известные алгоритмы, созданные в рамках различных подходов к конструированию методов оптимизации. Поскольку многие алгоритмы поиска экстремума многомерных функций основаны на редукции многомерной задачи к одной или семейству одномерных подзадач (см. раздел 5), то данная теория может быть применима и к анализу ряда методов многомерной оптимизации.

Рассмотрим одномерную задачу (4.6) для области поиска $Y = [a, b]$ и класс методов оптимизации (4.9), в котором $H_k(\Phi, \omega_k) = 1$ для любого $k \geq 1$, т.е. условие остановки отсутствует. В этом случае метод порождает бесконечную последовательность испытаний $\{y^k\} = y^1, y^2, \dots, y^k, \dots$, изучение свойств которой и будет предметом нашего интереса.

Определение 4.3. Алгоритм решения задачи (4.6) называется *характеристическим*, если, начиная с некоторого шага поиска $k_0 \geq 1$, выбор координаты y^{k+1} очередного испытания ($k \geq k_0$) заключается в выполнении следующих действий.

1) Задать набор

$$\Lambda_k = \{y_0, y_1, \dots, y_\tau\} \quad (4.46)$$

конечного числа $\tau + 1 = \tau(k) + 1$ точек области $Y = [a, b]$, полагая, что $a \in \Lambda_k, b \in \Lambda_k$, все координаты предшествующих испытаний $y^i \in \Lambda_k, 1 \leq i \leq k$, и множество Λ_k упорядочено (нижним индексом) по возрастанию координаты, т.е.

$$a = y_0 < y_1 < \dots < y_{\tau-1} < y_\tau = b. \quad (4.47)$$

2) Каждому интервалу $(y_{i-1}, y_i), 1 \leq i \leq \tau$, поставить в соответствие число $R(i)$, называемое характеристикой этого интервала.

3) Определить интервал (y_{t-1}, y_t) , которому соответствует максимальная характеристика $R(t)$, т.е.

$$R(t) = \max \{R(i) : 1 \leq i \leq \tau\} \quad (4.48)$$

4) Провести очередное испытание в точке

$$y^{k+1} = d(t) \in (y_{t-1}, y_t) \quad (4.49)$$

В соответствии с определением "характеристичность" алгоритма определяет структуру его решающего правила G_{k+1} через последовательность операций, представленных пунктами 1-4. Этим операциям можно дать следующую содержательную интерпретацию.

Для проведения нового испытания отрезок $[a, b]$ точками множества Λ_k разбивается на τ интервалов $(y_{i-1}, y_i), 1 \leq i \leq \tau$. Далее численно оценивается "перспективность" каждого интервала с помощью его характеристики и выбирается интервал, у которого характеристика наилучшая. Точка очередного испытания размещается внутри этого интервала в соответствии с правилом $d(\bullet)$.

Заметим, что множество (4.46) наряду с координатами испытаний может содержать точки, в которых испытания не проводились (например, в ряде информационно-статистических алгоритмов [9] такими точками являются концы отрезка). При этом *верхний индекс* координаты испытания соответствует *порядку* проведения испытаний в процессе поиска, а *нижний индекс* определяет *расположение* точки в упорядоченном наборе (4.47). Так, координата i -го испытания y^i в множестве Λ_k получит нижний индекс j , т.е. $y^i = y_j$, причем от шага к шагу номер j может меняться ($j = j(k)$).

Понятие характеристичности метода оптимизации впервые было введено В.А. Гришагиным [21] и позднее обобщено и распространено на другие классы задач и типы алгоритмов [15, 52, 53, 54].

В качестве иллюстрации приведем примеры известных алгоритмов глобальной оптимизации. В этих алгоритмах два первых испытания проводятся в точках $y^1 = a$ и $y^2 = b$, характеристическое правило вступает в действие,

начиная с $k=2$, при этом множество Λ_k ($k \geq 2$) состоит только из точек испытаний, т.е. $\Lambda_k = \{y^1, y^2, \dots, y^k\}$ и, следовательно, $\tau = k - 1$. Будем также использовать обозначение $z_j = f(y_j)$ для значений целевой функции в точках $y_j \in \Lambda_k$.

Метод последовательного сканирования (перебор).

Для этого метода характеристикой интервала является его длина, т.е.

$$R(i) = y_i - y_{i-1} \quad (4.50)$$

а точка очередного испытания выбирается в середине самого длинного интервала:

$$y^{k+1} = 0.5(y_{t-1} + y_t). \quad (4.51)$$

Метод Кушнера

Характеристика метода имеет вид

$$R(i) = - \frac{4(f_k^* - \delta_k - f(y_i))(f_k^* - \delta_k - f(y_{i-1}))}{y_i - y_{i-1}}, \quad (4.52)$$

а очередное испытание проводится в точке

$$y^{k+1} = y_{t-1} + \frac{(y_t - y_{t-1})(f_k^* - \delta_k - f(y_t))}{2(f_k^* - \delta_k) - f(y_t) - f(y_{t-1})}, \quad (4.53)$$

где $\delta_k > 0$ - параметр метода, в общем случае зависящий от номера шага поиска k , а f_k^* - наименьшее вычисленное значение функции.

Метод ломаных.

В данном методе, который был построен С.А.Пиявским [38] для оптимизации липшицевых функций $f \in Lip[a, b]$, характеристика

$$R(i) = 0.5m(y_i - y_{i-1}) - (z_i + z_{i-1})/2, \quad (4.54)$$

а точка очередного испытания выбирается согласно выражению

$$y^{k+1} = 0.5(y_t + y_{t-1}) - (z_t - z_{t-1})/(2m), \quad (4.55)$$

где $m > 0$ - параметр метода.

Информационно-статистический алгоритм глобального поиска (АГП).

Обсуждаемый метод предложен Р.Г.Стронгиным [7] как байесовский одношагово-оптимальный алгоритм и использует характеристику

$$R(i) = m(y_i - y_{i-1}) + \frac{(z_i - z_{i-1})^2}{m(y_i - y_{i-1})} - 2(z_i + z_{i-1}), \quad (4.56)$$

а точку нового испытания формирует согласно (4.42). Величина $m > 0$ вычисляется в соответствии с выражением

$$m = \begin{cases} rM, & M > 0 \\ 1, & M = 0 \end{cases} \quad (4.57)$$

где

$$M = \max_{1 \leq i \leq r} \frac{|z_i - z_{i-1}|}{y_i - y_{i-1}} \quad (4.58)$$

а $r > 1$ - параметр метода.

Контрольные вопросы и упражнения:

1. Выполните несколько первых итераций метода сканирования при минимизации функции $f(y)$ на отрезке $[a, b]$. Что произойдет, если последовательность испытаний будет бесконечной?
2. Постройте несколько первых точек последовательности поисковых испытаний АГП при решении задачи минимизации линейной функции $f(y) = y$ на отрезке $[0, 1]$, принимая параметр $r=2$. Попробуйте установить аналитическую закономерность размещения точек испытаний в этой задаче.
3. Выполните задание 2 для метода ломаных, применяя в качестве параметра метода m оценку (4.44) при $r=2$. Попробуйте установить связь между последовательностями испытаний метода ломаных и АГП.

После конкретных примеров представим общий теоретический результат о связи принципа одношаговой оптимальности со свойством характеристичности.

Теорема 4.1. *Одношагово-оптимальный (минимаксный или байесовский) алгоритм является характеристическим.*

ДОКАЗАТЕЛЬСТВО. Перепишем условие одношаговой оптимальности (4.32) в виде

$$W_k(y^k) = \min_{\tilde{y}^k \in [a, b]} W_k(\tilde{y}^k) = \min_{1 \leq i \leq r} \min_{\tilde{y}^k \in [y_{i-1}, y_i]} W_k(\tilde{y}^k). \quad (4.59)$$

Отсюда следует, что в качестве характеристики интервала (y_{i-1}, y_i) можно взять величину

$$R(i) = - \min_{\tilde{y}^k \in [y_{i-1}, y_i]} W_k(\tilde{y}^k), \quad (4.60)$$

а

$$d(t) = \arg \min_{\tilde{y}^k \in [y_{i-1}, y_i]} W_k(\tilde{y}^k).$$

Теорема доказана.

Теорема 4.2. *Пусть точка y^* является предельной точкой (точкой накопления) последовательности поисковых испытаний $\{y^k\}$, порождаемой характеристическим алгоритмом при решении задачи (4.6) на отрезке $[a, b]$, причем $y^* \neq a$ и $y^* \neq b$. Предположим, что характеристики $R(i)$ и правила выбора точки очередного испытания $d(t)$ удовлетворяют следующим требованиям:*

а) если при $k \rightarrow \infty$ точка $\bar{y} \in [y_{i(k)-1}, y_{i(k)}]$ и $y_{i(k)-1} \rightarrow \bar{y}$, $y_{i(k)} \rightarrow \bar{y}$, тогда

$$R(i(k)) \rightarrow -\mu f(\bar{y}) + c; \quad (4.61)$$

b) в случае, когда, начиная с некоторого шага поиска, интервал $(y_{i-1}, y_i), i = i(k)$ не содержит точек поисковых испытаний, т.е. существует номер $\tilde{k} \geq 1$ такой, что для всех $k \geq \tilde{k}$

$$(y_{i-1}, y_i) \cap \bigcap \{y^k\} = \emptyset, \quad (4.62)$$

для характеристики интервала справедливо

$$\lim_{k \rightarrow \infty} R(i) > -\mu \min\{f(y_{i-1}), f(y_i)\} + c, \quad (4.63)$$

c) при выборе очередного испытания имеет место соотношение

$$\max\{y^{k+1} - y_{t-1}, y_t - y^{k+1}\} \leq \nu(y_t - y_{t-1}), \quad (4.64)$$

где μ, c, ν - некоторые константы, причем $\mu \geq 0, 0 < \nu < 1$.

Тогда последовательность $\{y^k\}$ содержит две подпоследовательности, одна из которых сходится к y^* слева, а другая справа.

ДОКАЗАТЕЛЬСТВО. Рассмотрим вначале случай, когда $y^* \notin \{y^k\}$. Обозначим через $p = p(k), k \geq 1$, номер интервала (y_{p-1}, y_p) , содержащего на k -м шаге поиска предельную точку y^* . Очевидно, что для $k = 1$ $[y_{p-1}, y_p] = [a, b]$. После попадания в интервал (y_{p-1}, y_p) точки очередного испытания y^{k+1} (в этом случае $p = t$) для нового интервала $(y_{p(k+1)-1}, y_{p(k+1)})$, содержащего y^* , согласно (4.64) справедлива оценка

$$y_{p(k+1)-1} - y_{p(k+1)} \leq \nu(y_{p(k)-1} - y_{p(k)})$$

Но тогда после попадания с начала поиска s испытаний в интервал с точкой y^* его длина будет удовлетворять неравенству

$$y_{p-1} - y_p \leq \nu^s(b-a). \quad (4.65)$$

Поскольку точка y^* предельная, после образования на некотором шаге интервала (y_{p-1}, y_p) в него попадет бесконечное число испытаний, поэтому из (4.65) следует, что

$$\lim_{k \rightarrow \infty} (y_{p(k)-1} - y_{p(k)}) = 0. \quad (4.66)$$

На основании (4.66) в качестве искомым подпоследовательностей мы можем взять последовательность $\{y_{p(k)-1}\}$ левых и последовательность $\{y_{p(k)}\}$ правых концов интервалов, содержащих y^* .

Пусть теперь найдется номер $q \geq 1$ такой, что $y^q = y^*$. Тогда при любом $k \geq q$ существует номер $j = j(k), 0 \leq j \leq \tau$, для которого $y_j = y^*$. Допустим, что имеет место односторонняя сходимости к y^* , например, слева. Тогда найдется номер $\tilde{k} \geq q$ такой, что при $k \geq \tilde{k}$ испытания в интервал (y_j, y_{j+1}) попадать не будут.

Из (4.62), (4.63) для интервала (y_j, y_{j+1}) и как следствие соотношения (4.61) для интервала (y_{j-1}, y_j) мы получаем, что

$$\lim_{k \rightarrow \infty} R(j+1) > -\mu \min(f(y_j), f(y_{j+1})) + c \geq -\mu f(y^*) + c$$

$$\lim_{k \rightarrow \infty} R(j) = -\mu f(y^*) + c$$

откуда, начиная с некоторого шага поиска, будет следовать выполнимость неравенства

$$R(j+1) > R(j) \quad (4.67)$$

Однако вследствие решающего правила (4.46)-(4.49) соотношение (4.67) противоречит невозможности проведения испытаний в интервале (y_j, y_{j+1}) . Данное противоречие завершает доказательство.

Следствие 4.2.1. В вычислительную схему характеристического алгоритма, удовлетворяющую условиям Теоремы 4.2, можно ввести условие остановки вида

$$y_t - y_{t-1} \leq \varepsilon, \quad (4.68)$$

где t из (4.48), а $\varepsilon > 0$ - заданная точность поиска (по координате), т.е. прекращать вычисления, когда длина интервала с максимальной характеристикой станет меньше заданной точности ε . Тогда процесс поиска будет остановлен через конечное число шагов.

ДОКАЗАТЕЛЬСТВО. Вначале укажем, что на конечном отрезке $[a, b]$ последовательность испытаний всегда будет иметь хотя бы одну предельную точку y^* . Обозначим через $p = p(k), k \geq 1$, номер интервала, содержащего точку y^* на k -м шаге поиска. Т.к. данная точка предельная, то в интервал (y_{p-1}, y_p) попадет бесконечное число испытаний и для него будет иметь место соотношение (4.66), из которого следует, что условие (4.68) неизбежно выполнится на некотором шаге поиска. Заметим, что если точка y^* не является внутренней, для справедливости (4.68) достаточно односторонней сходимости.

Проверим выполнимость условий Теоремы 4.2 для рассмотренных примеров характеристических алгоритмов.

Метод последовательного сканирования.

Если интервал (y_{i-1}, y_i) стягивать в точку, характеристика метода (4.39) стремится к нулю. Поэтому в (4.61) в качестве констант μ и c можно взять $\mu = c = 0$. Если же в интервал (y_{i-1}, y_i) , начиная с некоторого шага поиска, испытания попадать не будут, то его длина (совпадающая с характеристикой), будет оставаться положительной, что обеспечит выполнимость условия (4.63) для выбранных μ и c . Что касается неравенства (4.64), то оно очевидно выполняется при $\nu = 0.5$.

Метод Кушнера

Предположим, что функция $f(y)$ ограничена на отрезке $[a, b]$ конечными величинами f_{\min} и f_{\max} , т.е. $f_{\min} \leq f(y) \leq f_{\max}, y \in [a, b]$.

Если для любого шага k параметр метода $\delta_k > \delta > 0$, то $f_k^* - f(y_i) - \delta_k < -\delta < 0$ и $f_k^* - f(y_{i-1}) - \delta_k < -\delta < 0$, то

$$R(i(k)) = -\frac{4(f_k^* - \delta_k - f(y_i))(f_k^* - \delta_k - f(y_{i-1}))}{y_i - y_{i-1}} \rightarrow -\infty,$$

когда длина интервала (y_{i-1}, y_i) стремится к нулю. Поэтому (4.61) выполняется при $\mu = 0$ и формальном значении $c = -\infty$ (Доказательство Теоремы 4.2 несколько не изменится, если в правой части (4.61) предел будет равен $-\infty$).

Поскольку для любого интервала с ненулевой длиной его характеристика будет иметь конечное значение, (4.63) также справедливо.

Рассмотрим разности

$$y^{k+1} - y_{t-1} = \beta_1(y_t - y_{t-1}), \quad y_t - y^{k+1} = \beta_2(y_t - y_{t-1}), \quad \text{где}$$

$$\beta_1 = \frac{f(y_{t-1}) - f_k^* + \delta_k}{f(y_{t-1}) + f(y_t) - 2(f_k^* - \delta_k)}, \quad \beta_2 = \frac{f(y_t) - f_k^* + \delta_k}{f(y_{t-1}) + f(y_t) - 2(f_k^* - \delta_k)}.$$

Введем вспомогательную функцию $\varphi(x, \alpha) = \frac{x + \alpha}{x + 2\alpha}$. Так как $\min\{f(y_{t-1}), f(y_t)\} \geq f_k^*$, то $\beta_1 \leq \varphi(f(y_{t-1}) - f_k^*, \delta_k)$, $\beta_2 \leq \varphi(f(y_t) - f_k^*, \delta_k)$.

При положительных x функция $\varphi(x, \alpha)$ убывает по α , ибо $\varphi'_\alpha(x, \alpha) = \frac{-x}{(x + 2\alpha)^2} < 0$, поэтому для $\delta_k > \delta > 0$ выполняется

$$\beta_1 \leq \varphi(f(y_{t-1}) - f_k^*, \delta), \quad \beta_2 \leq \varphi(f(y_t) - f_k^*, \delta).$$

В свою очередь функция $\varphi(x, \alpha)$ при положительных α возрастает по x , так как $\varphi'_x(x, \alpha) = \frac{\alpha}{(x + 2\alpha)^2} > 0$, поэтому

$$\max\{\varphi(f(y_{t-1}) - f_k^*, \delta), \varphi(f(y_t) - f_k^*, \delta)\} \leq \varphi(f_{\max} - f_{\min}, \alpha).$$

Последнее неравенство означает, что

$$\max\{\beta_1, \beta_2\} \leq \varphi(f_{\max} - f_{\min}, \alpha) = \frac{f_{\max} - f_{\min} + \delta}{f_{\max} - f_{\min} + 2\delta}.$$

Именно последнюю величину можно выбрать в качестве константы ν для соотношения (4.64), ибо очевидно, что $0 < \nu < 1$.

Метод ломаных

Для рассмотрения условий теоремы сделаем предположение, что минимизируемая функция $f(y)$ удовлетворяет условию Липшица с константой $L > 0$, т.е.

$$|f(y') - f(y'')| \leq L|y' - y''|, \quad y', y'' \in [a, b] \quad (4.69)$$

и, кроме того, параметр метода $m > L$.

В силу липшицевости функция $f(y)$ непрерывна, и, следовательно, при стягивании интервала (y_{i-1}, y_i) к точке \bar{y} характеристика интервала будет стремиться к величине $-f(\bar{y})$, т.е. для выполнимости (4.61) можно положить $\mu = 1$ и $c = 0$. Проверим теперь для данных μ и c справедливость неравенства (4.50) для интервала (x_{i-1}, x_i) , удовлетворяющего (4.63). Воспользуемся простым соотношением

$$\min\{z_{i-1}, z_i\} = \frac{1}{2}(z_{i-1} + z_i - |z_{i-1} - z_i|) \quad (4.70)$$

и оценим характеристику (4.54) метода, воспользовавшись условием Липшица и учитывая, что длина интервала остается, начиная с некоторого шага поиска, положительной и неизменной:

$$0.5m(y_i - y_{i-1}) - (z_i + z_{i-1})/2 > 0.5L(y_i - y_{i-1}) - (z_i + z_{i-1})/2 \geq 0.5(z_{i-1} + z_i - |z_{i-1} - z_i|) =$$

$$= -\min\{z_{i-1}, z_i\}$$

Полученные соотношения устанавливают справедливость (4.63).

Для определения величины ν в (4.64) оценим величину

$$y_i - y^{k+1} = 0.5(y_i - y_{i-1}) + (z_i - z_{i-1})/(2m) \leq 0.5(y_i - y_{i-1}) + L(y_i - y_{i-1})/(2m) =$$

$$= 0.5(1 + L/m)(y_i - y_{i-1})$$

Аналогичная оценка имеет место и для интервала (y_{i-1}, y^{k+1}) , поэтому можно взять $\nu = 0.5(1 + L/m)$. Очевидно, что $\nu > 0$ и, поскольку $m > L$, то $\nu < 1$.

Информационно-статистический алгоритм глобального поиска

Предположив липшицевость (4.69) целевой функции, нетрудно показать, что данный метод также удовлетворяет условиям теоремы 4.2 с $\mu=4$, $c=0$ и $\nu = 0.5(1 + 1/r)$.

Вследствие липшицевости при стягивании интервала (y_{i-1}, y_i) к точке \bar{y} его характеристика будет стремиться к величине $-4f(\bar{y})$, т.е. для выполнимости (4.61) можно положить $\mu=4$ и $c=0$.

Для проверки (4.63) при условии (4.62) рассмотрим два случая. Пусть сначала для интервала (y_{i-1}, y_i) справедливо $z_{i-1} = z_i$. Тогда характеристика

$$R(i) = m(y_i - y_{i-1}) - 2(z_i + z_{i-1}) > -2(z_i + z_{i-1}) = -4\min\{z_{i-1}, z_i\},$$

поскольку длина интервала (y_{i-1}, y_i) , начиная с некоторого шага поиска k_Δ , перестает изменяться, т.е. существует константа $\Delta > 0$ такая, что при $k > k_\Delta$ имеет место $y_i - y_{i-1} > \Delta > 0$.

Предположим теперь, что $z_{i-1} \neq z_i$. Представим характеристику (4.56) в виде

$$R(i) = |z_i - z_{i-1}|(\beta + \frac{1}{\beta}) - 2(z_i + z_{i-1}), \quad (4.71)$$

где вследствие (4.57), (4.58) $0 < \beta = \frac{m(y_i - y_{i-1})}{|z_i - z_{i-1}|} < 1$. В этом случае величина

$$\beta + \frac{1}{\beta} > 2, \text{ поэтому}$$

$$R(i) > 2|z_i - z_{i-1}| - 2(z_i + z_{i-1}) = -4\min\{z_i, z_{i-1}\}.$$

Что касается условия (4.64), то в соответствии с (4.58) справедлива оценка

$$y_i - y^{k+1} = 0.5(y_i - y_{i-1}) + (z_i - z_{i-1})/(2m) \leq 0.5(y_i - y_{i-1}) + M(y_i - y_{i-1})/(2m) =$$

$$= 0.5(1 + 1/r)(y_i - y_{i-1})$$

Длину интервала (y_{i-1}, y^{k+1}) можно оценить аналогичным образом, поэтому в качестве ν можно выбрать величину $0.5(1 + 1/r)$, которая очевидным образом удовлетворяет условию $0 < \nu < 1$, т.к. $r > 1$.

Теорема 4.3. Если в условиях теоремы 4.2 $\mu = 0$, тогда любая точка области поиска является предельной точкой последовательности поисковых испытаний.

ДОКАЗАТЕЛЬСТВО. Предположим, что некоторая точка y' не является предельной для последовательности поисковых испытаний, порождаемой алгоритмом. Это означает, что начиная с некоторого шага поиска, испытания в интервал (y_{q-1}, y_q) , $q = q(k)$, содержащий эту точку, попадать не будут, и тогда согласно (4.63) характеристика $R(q) > 0$. Но последовательность испытаний, будучи ограниченной пределами отрезка $[a, b]$, содержит хотя бы одну сходящуюся подпоследовательность. Для интервала (y_{p-1}, y_p) , $p = p(k)$, содержащего предельную точку данной подпоследовательности, вследствие двусторонней сходимости справедливо соотношение (4.61), т.е. его длина должна стремиться к нулю. Это означает, что на некотором шаге характеристика $R(p)$ станет меньше характеристики $R(q)$, что в соответствии с правилом (4.48) основной схемы характеристического алгоритма противоречит исходному предположению.

Теорема доказана.

Данная теорема устанавливает условия так называемой "всюду плотной" сходимости, когда метод сходится ко всем точкам области поиска, в том числе, разумеется, и к точкам глобального минимума. Указанный тип сходимости адекватен таким классам задач, для которых невозможно построить оценки экстремума по конечному числу испытаний (например, для класса непрерывных функций), и в этом случае обеспечить сходимость к глобально-оптимальному решению можно лишь за счет свойства всюду плотной сходимости.

Среди примеров, которые мы рассмотрели, подобной сходимостью обладают методы перебора и метод Кушнера, а среди других известных алгоритмов – методы [15].

Для иллюстрации поведения методов данного типа продемонстрируем динамику поиска метода перебора. Под графиком минимизируемой функции штрихами отображены координаты проведенных испытаний.

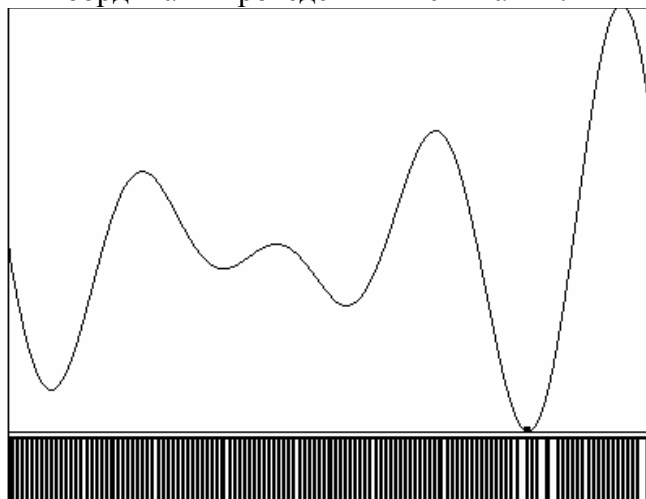


Рис.4.2. Размещение испытаний в методе равномерного перебора. *Всюду плотная сходимость*

Условия всюду плотной сходимости обеспечивают достаточные условия сходимости к глобально-оптимальному решению задачи оптимизации. Однако характер такого рода сходимости требует дополнительных способов

исследования эффективности распределения точек поисковой последовательности. Один из таких подходов основан на получении оценок относительной плотности размещения испытаний в различных подобластях области поиска в сравнении с плотностью точек в окрестности глобального минимума [15].

Метод аналитического построения оценок относительной плотности размещения испытаний будет рассмотрен позднее, в разделе 4.6 настоящей главы. Здесь же мы ограничимся тем, что приведем вид этих оценок на примере всюду плотно сходящегося метода Х.Кушнера. Оценка получена на классе непрерывных кусочно-линейных функций.

Пусть $\alpha_k(y, y^*)$ — относительная плотность испытаний в точке y по отношению к точке решения y^* . Можно доказать, что для каждого y , начиная с некоторого k , относительная плотность $\alpha_k(y, y^*)$ совершает колебания вокруг значения

$$\tilde{\alpha}(y, y^*) = \left(\frac{1}{1 + ((f(y) - f^*)/\delta)} \right)^2 \quad (4.72)$$

не покидая сколь угодно малую окрестность интервала $[0.5 \tilde{\alpha}(y, y^*); 2 \tilde{\alpha}(y, y^*)]$, причем при значениях, больших чем (4.72), относительная плотность только убывает, а при меньших значениях — только возрастает. Следует обратить внимание на тот факт, что глобальным минимумам задачи соответствуют наибольшие значения относительной плотности (4.72), а локальным минимумам функции $f(y)$ — локальные максимумы относительной плотности. Из (4.72) следует, что чем меньше значение δ , тем более выражена наибольшая концентрация испытаний в окрестностях глобальных минимумов задачи. Соответствующая иллюстрация приведена на рис. 4.6 в разделе 4.6.

На рис.4.3 приведен пример размещения испытаний всюду плотно сходящегося метода Х.Кушнера с автоматической настройкой параметра δ_k по алгоритму, подобному (4.43), при $\delta_0 = 1, \beta = 0.1, \varepsilon^f = 0.01$. В процессе поиска параметр настраивался следующим образом: $\delta_9 = 0.5, \delta_{20} = 0.2, \delta_{34} = 0.05, \delta_{49} = 0.02, \delta_{52} = 0.01$. Иллюстрация ясно показывает существенно неравномерный характер размещения испытаний в этом методе (с их концентрацией в окрестности решения), несмотря на всюду плотный тип сходимости.

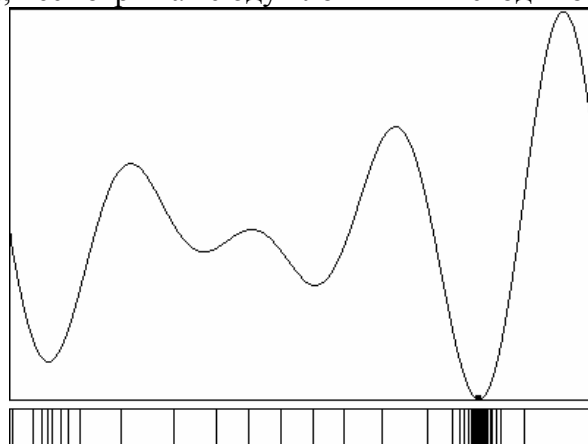


Рис.4.3. Характер концентрации испытаний в методе Х.Кушнера. Сходимость в пределе всюду плотная, но неравномерная

Другой тип поведения характеристических алгоритмов устанавливает

Теорема 4.4. Пусть в условиях теоремы 4.2 $\mu > 0$ и y^* - предельная точка последовательности поисковых испытаний. Тогда

1) $f(y^k) \geq f(y^*), k \geq 1$;

2) если существует еще одна предельная точка $y^{**} \neq y^*$, то $f(y^{**}) = f(y^*)$;

3) если в области поиска функция $f(y)$ имеет конечное число локальных минимумов, то y^* является точкой локального минимума целевой функции в области Y .

ДОКАЗАТЕЛЬСТВО. Докажем вначале второе утверждение. Пусть $f(y^{**}) \neq f(y^*)$ и для определенности $f(y^{**}) > f(y^*)$. Обозначим через $p=p(k)$ и $q=q(k)$ номера интервалов (y_{p-1}, y_p) и (y_{q-1}, y_q) , и y^{**} соответственно. Вследствие Теоремы 4.2 и соотношения (4.61)

$$\lim_{k \rightarrow \infty} R(p(k)) = -\mu f(y^*) + c,$$

$$\lim_{k \rightarrow \infty} R(q(k)) = -\mu f(y^{**}) + c,$$

Это означает, что начиная с некоторого шага поиска, начнет выполняться неравенство

$$R(p) > R(q),$$

которое приведет к тому, что согласно правилу (4.48) испытания в интервал (y_{q-1}, y_q) попадать не будут, что противоречит предположению о том, что y^{**} - предельная точка.

Перейдем к доказательству первого утверждения.

Предположим, что на некотором шаге с номером $j \geq 1$ в точке y^j получено значение $f(y^j) < f(y^*)$. Обозначим через $s=s(k)$, $k \geq j$, номер точки y^j в множестве (4.46), т.е. $y_s = y^j$. Покажем, что точка y^j - тоже предельная. Если это не так, то характеристика интервала (y_{s-1}, y_s) (если $y^j = a$, то надо рассмотреть интервал (y_s, y_{s+1})) согласно (4.63) удовлетворяет неравенству

$$\lim_{k \rightarrow \infty} R(s) > -\mu f(y^j) + c.$$

Поскольку для интервала (y_{p-1}, y_p) , содержащего на k -м шаге поиска точку y^* , вследствие двусторонней сходимости и (4.61) справедливо

$$\lim_{k \rightarrow \infty} R(p(k)) = -\mu f(y^*) + c,$$

то, начиная с некоторого шага поиска, характеристика $R(s)$ интервала (y_{s-1}, y_s) станет больше характеристики $R(p)$ и в интервал (y_{p-1}, y_p) испытания из-за правила (4.48) попадать больше не смогут, а это противоречит предположению о том, что точка — предельная. Но после этого наступает немедленное противоречие со вторым утверждением теоремы, и это противоречие доказывает наше первое утверждение.

Завершим теорему доказательством локальной оптимальности точки y^* . Предположим обратное, а именно, что y^* не является точкой локального оптимума. Тогда существует непустая окрестность $\Theta \subseteq [a, b]$ точки y^* , в которой функция $f(y)$ строго монотонна, то есть либо слева, либо справа от y^* имеет место неравенство $f(y) < f(y^*)$, $y \in \Theta$. Так как сходимость к точке y^* двусторонняя, в последовательности испытаний обязательно появится точка $y^j \in \Theta$ такая, что $f(y^j) < f(y^*)$, что противоречит первому утверждению и завершает доказательство теоремы.

Итак, если в условиях (4.61), (4.63) константа μ положительна (а это означает, что метод *учитывает* в асимптотике информацию о функции), то поведение алгоритма становится более целенаправленным: он ищет только такие точки, которые обладают свойством локальной минимальности. Более того, утверждение 2 теоремы говорит, что метод не может сходиться к разновысоким локальным минимумам.

Из числа рассмотренных нами методов таким свойством обладают метод ломаных и АГП, а из других известных алгоритмов – весь спектр информационно-статистических алгоритмов [7, 9, 42, 53], а также методы [52].

Для примера проиллюстрируем работу АГП. Как и ранее, штрихами под графиком функции отмечаются координаты испытаний метода.

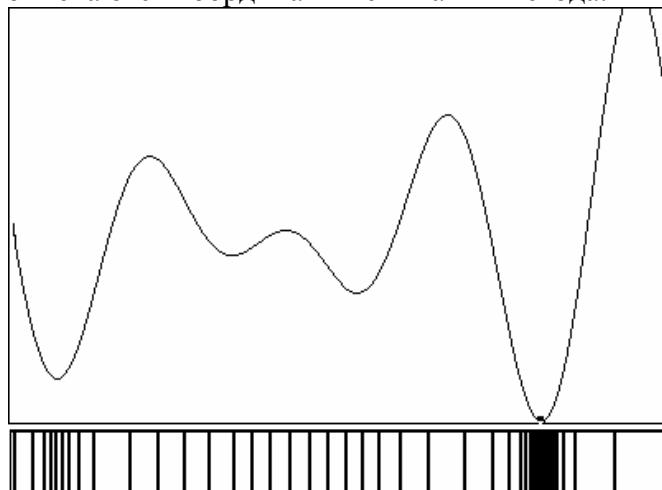


Рис.4.4. Размещение испытаний в методе АГП

Приведенный рисунок показывает, что в той части области поиска, где нет глобального минимума, метод строит *редкую сетку* испытаний, а сгущение, обусловленное сходимостью, наблюдается лишь в окрестности глобального минимума.

Теорема 4.4 не гарантирует сходимость к глобальному минимуму исследуемой задачи оптимизации. Такие гарантии (достаточные условия сходимости) дает

Теорема 4.5. Пусть функция $f(y)$ удовлетворяет на отрезке $[a, b]$ условию Липшица (4.69). Тогда любая точка y^* глобального минимума функции $f(y)$ на данном отрезке является предельной точкой последовательности поисковых испытаний, порождаемой характеристическим алгоритмом при решении задачи

(4.6), если выполняются условия (4.61) и (4.64) теоремы 4.2, а для интервала (y_{i-1}, y_i) со свойством (4.62) справедливо неравенство

$$\lim_{k \rightarrow \infty} R(i) > \frac{\mu}{2}(L(y_i - y_{i-1}) - f(y_{i-1}) - f(y_i)) + c, \quad (4.73)$$

где константы μ, c и ν удовлетворяют условиям теоремы 4.2.

ДОКАЗАТЕЛЬСТВО. Предположим, что y^* не является предельной точкой поисковых испытаний и пусть $i=p(k)$ есть номер интервала (y_{i-1}, y_i) , содержащего эту точку на k -м шаге поиска, то есть $y_{i-1} \leq y^* \leq y_i$. Тогда, начиная с некоторого шага поиска, испытания в данный интервал попадать не будут. Находясь в интервале (y_{i-1}, y_i) , точка y^* может быть записана в виде $y^* = \lambda y_{i-1} + (1 - \lambda)y_i, 0 \leq \lambda \leq 1$. Тогда вследствие (4.69) имеем

$$\begin{aligned} f(y_{i-1}) - f(y^*) &\leq L(y^* - y_{i-1}) = L(1 - \lambda)(y_i - y_{i-1}), \\ f(y_i) - f(y^*) &\leq L(y^* - y_i) = L\lambda(y_i - y_{i-1}), \end{aligned}$$

откуда

$$f(y_{i-1}) + f(y_i) \leq 2f(y^*) + L(y_i - y_{i-1}).$$

Последнее неравенство в сочетании с (4.73) позволяет получить оценку

$$\lim_{k \rightarrow \infty} R(i) > -\mu f(y^*) + c.$$

В то же время последовательность испытаний имеет на отрезке $[a, b]$ хотя бы одну предельную точку \bar{y} , для которой необходимо выполнено $f(\bar{y}) \geq f(y^*)$ и, кроме того, согласно Теореме 4.2, для интервала $(y_{p-1}, y_p), p=p(k)$, содержащего \bar{y} , имеет место соотношение

$$\lim_{k \rightarrow \infty} R(p) = -\mu f(\bar{y}) + c$$

Но тогда, начиная с некоторого шага поиска, получим неравенство $R(i) > R(p)$, которое не допускает размещение испытаний в интервале (y_{p-1}, y_p) и противоречит тому, что точка \bar{y} - предельная.

Доказательство завершено.

Следствие 4.5.1. При $\mu = 0$ (4.73) совпадает с (4.63) и теорема 4.5 становится идентичной теореме 4.3; при этом липшицевость целевой функции необязательна.

Следствие 4.5.2. Если в условиях теоремы $\mu > 0$, то характеристический алгоритм сходится ко всем точкам глобального минимума и только к ним.

ДОКАЗАТЕЛЬСТВО. Действительно, сходимость ко всем точкам глобального минимума – это результат самой теоремы, а невозможность сходимости к точкам, отличным от глобально оптимальных, - следствие утверждения 2 теоремы 4.4.

Заметим, что для метода ломаных (4.73) очевидно выполняется при $m > L$.

Для характеристики АГП справедливо очевидное неравенство

$$R(i) \geq m(y_i - y_{i-1}) - 2(z_i + z_{i-1}) = m(y_i - y_{i-1}) - \frac{\mu}{2}(z_i + z_{i-1}),$$

откуда следует выполнимость (4.73) при $m > 2L$.

4.5. Анализ сходимости многомерных методов многоэкстремальной оптимизации. T -представимые алгоритмы и их свойства

Характеристическая теория нашла свое развитие в ряде других теоретических обобщений, из числа которых остановим свое внимание на теории типичной представимости (сокращенно T -представимости), предложенной С.Ю. Городецким [15] и ориентированной в первую очередь, на анализ сходимости многомерных алгоритмов.

Основная идея этого подхода состоит в том, чтобы свести анализ сходимости (всюду плотная, сходимость только к подмножеству локальных минимумов и т.д.) к анализу свойств функционала $W_k(y)$, определяющего правило размещения очередного испытания. Причем для облегчения такого анализа используется существенно более простой «предельный» функционал $W^*(y)$, вычисляемый по исходному.

Прежде чем привести более точные определения поясним на простых примерах смысл «предельного» функционала. Рассмотрим выражение $W_{k+1}(y)$ из (4.34) для информационно–статистического алгоритма. Введем бесконечную последовательность точек y^i , всюду плотно в пределе покрывающую область определения $D=[a;b]$ целевой функции. Если $W_k(y)$ строить по результатам измерений целевой функции в точках этой искусственно введенной последовательности, то из вида (4.34) очевидно, что при $k \rightarrow \infty \forall y \in D$ будет существовать поточечный предел

$$W_k(y) \rightarrow W^*(y) = -4f(y).$$

В данном случае это и будет предельный функционал. Можно заметить (в последующем это окажется важным), что его значения связаны со значениями целевой функции строго монотонным преобразованием.

Заметим, что слово «функционал», употребляемое в этом разделе, подчеркивает, что $W_k(\cdot)$ и $W_k^*(\cdot)$ зависят от функций f из некоторого класса Φ .

Если рассмотреть метод Х.Кушнера, то при тех же условиях, как легко видеть из (4.36), при $k \rightarrow \infty \forall y \in D$ будет существовать поточечный предел

$$W_k(y) \rightarrow W^*(y) = const = -\infty.$$

Таким образом, для метода Х.Кушнера предельный функционал является константой, что принципиально отличает его от информационно–статистического метода.

Ниже будет показано, что вид предельного функционала существенно влияет на тип сходимости.

4.5.1. Описание класса задач

Перейдем теперь к формальному изложению теории. Будем рассматривать задачу оптимизации с целевой функцией $f(y)$. В задаче могут присутствовать ограничения–неравенства $g(y) \leq 0$, где $g(y) = (g_1(y), \dots, g_m(y))$. Все функции считаются заданными в многомерном параллелепипеде D . Заметим, что можно рассмотреть случай, когда функциональных ограничений нет, и даже случай, когда они присутствуют, но зато нет целевой функции (при этом задача оптимизации превратится в задачу поиска решения системы неравенств). Для

единообразного описания этих классов задач введем функцию Q (в общем случае, она будет вектор–функцией), компонентами которой будут являться все функции, входящие в постановку задачи. Например, в наиболее полном случае, $Q = (f, g_1, \dots, g_m)$.

Пусть Y^* — множество решений рассматриваемой задачи (например, множество глобальных минимумов, если рассматривается задача оптимизации). В общем случае по измерению значения $Q(y)$ нельзя определить принадлежность точки y к Y^* . Однако существуют задачи, в которых это возможно. Например, задачи оптимизации с известным глобально–оптимальным значением f^* целевой функции. Для таких задач по значению $Q(y)$ можно определить принадлежность y к Y^* . В таких случаях будем говорить, что у задачи существует индикатор множества решений $I(y)$. Определим его следующим образом

$$I(y) = \begin{cases} 0, & y \notin Y^* \\ 1, & y \in Y^*. \end{cases} \quad (4.74)$$

Дополнительным примером задач с известным индикатором могут служить задачи решения совместных систем уравнений.

4.5.2. Класс Т–представимых алгоритмов, классификация, примеры

Рассмотрим следующий класс алгоритмов поиска решения.

Для точек поисковых испытаний, как и ранее, сохраним обозначения $y^i, i = 1, 2, \dots$. Множество первых k точек обозначим как Y_k , а бесконечную последовательность испытаний – через Y_∞ . Испытания будут состоять в измерении значений функции Q (хотя можно рассматривать и измерения более общего вида). Набор вычисленных значений Q вместе с координатами точек испытаний составляют, как и раньше, поисковую информацию $\omega_k = \omega_k(Q) = \omega(Q, Y_k)$. В последующем будем использовать все три приведенные формы записи. Последняя форма удобна тем, что позволяет явно указать, в каких точках проводились измерения.

Примем, что правило выбора точек новых испытаний y^k , начиная с некоторого шага k_0 (т.е. при $k > k_0$), определяется соотношениями

$$y^{k+1} = \arg \max \{ W_k(\omega(Q, Y_k), y) : y \in D \}, \quad Y_{k+1} = Y_k \cup \{ y^{k+1} \}. \quad (4.75)$$

Функционал $W_k(\omega(Q, Y_k), y)$, заданный на функциях Q из некоторого класса Φ , может быть построен, например, на основе принципов одношаговой оптимальности с помощью соотношений (4.31) или (4.33).

 **Замечание.** В виде (4.75) представимо любое решающее правило (4.14).

ДОКАЗАТЕЛЬСТВО. Справедливость этого замечания становится очевидной, если положить

$$W_k(\omega_k(Q), y) = -\|G_{k+1}(\Phi, \omega_k(Q)) - y\|.$$

Введем ряд дополнительных обозначений. Произвольную бесконечную последовательность различных точек $u^i, i = 1, 2, \dots$ из области поиска D обозначим через U_∞ , а множество ее первых k членов — через U_k , т.е.

$$U_\infty = \{u^i \in D : i = 1, 2, \dots\}, \quad U_k = \{u^i \in D : 1 \leq i \leq k\}. \quad (4.76)$$

Результаты измерения функций задачи на начальных участках произвольных последовательностей точек будем записывать как $\omega(Q, U_k)$

Определение 4.4. Процедуру поиска назовем типично-представимой (Т-представимой) на классе Φ , если, начиная с некоторого шага поиска, правило выбора точки очередного испытания представимо в виде (4.75) с функционалом $W_k(\omega(Q, U_k), y)$, для которого существует другой, предельный функционал $W^*(y, Q(y))$, удовлетворяющий для $Q \in \Phi$ следующим условиям.

1. Для любой точки $y \in D$ и последовательности U_∞ из (4.76), содержащей y среди своих предельных точек, выполняется

$$\lim_{k \rightarrow \infty} W_k(\omega(Q, U_k), y) = W^*(y, Q(y)). \quad (4.77)$$

2. Для задач с неизвестным индикатором $I(y)$ множества решений – при всех $y \in D$, а для задач с известным индикатором – только для $y \in D \setminus Y^*$ на всякой сходящейся к y подпоследовательности точек $u^{i_k} \in U_\infty, k=1,2,\dots$ имеет место

$$\lim_{k \rightarrow \infty} W_{i_k}(\omega(Q, U_{i_k}), u^{i_k}) = W^*(y, Q(y)). \quad (4.78)$$

3. Для любого $u \in U_k$

$$W_k(\omega(Q, U_k), u) = W^*(u, Q(u)). \quad (4.79)$$

Заметим, что произвольная последовательность U_∞ из (4.76) использована в определении 4.4 для того, чтобы можно было изучать свойства функционала W_k , не зная порождаемой им, в силу (4.75), последовательности поисковых испытаний.

Функционал W_k из (4.75) назовем *характеристическим*, а W^* — *предельным*.

К классу Т-представимых относятся многие известные алгоритмы, например, рассмотренные в разделах 4.3–4.4 одношагово–оптимальные методы: информационно–статистический метод Р.Г.Стронгина, метод ломаных С.А.Пиявского, байесовский метод Х.Кушнера, а также методы, разработанные позднее, например, многомерный метод условной оптимизации на основе адаптивных вероятностных моделей, рассмотренный в пункте 4.3.5.

Предельный функционал W^* несет в себе важную информацию о поведении поисковой последовательности, порождаемой правилом выбора точек испытаний (4.75). Можно дать неполную классификацию Т-представимых процедур по свойствам предельного функционала, выделив основные важные подклассы. В последующем они будут увязаны с определенными типами сходимости

Процедурами *константного типа* (Const–типа) назовем Т-представимые процедуры, у которых предельный функционал является константой

$$W^*(y, Q(y)) = const \geq -\infty.$$

К δ -типу отнесем Т-представимые процедуры с

$$W^*(y, Q(y)) = \eta(I(y)),$$

где $I(y)$ — индикатор множества решений Y^* и $\eta(0) < \eta(1)$.

Например, если предельный функционал равен 0 вне множества решений и 1 на Y^* , то Т-представимая процедура относится к δ -типу.

Подкласс процедур *непрерывного типа* (C -типа¹) характеризуется условием

$$W^*(y, Q(y)) = \mu(Q(y)),$$

где функция μ непрерывна.

К *монотонному типу* (M -типу) отнесем методы непрерывного типа (C -типа) со скалярной функцией $Q(y) = f(y)$, в которых $\mu(z)$ — монотонно невозрастающая функция z , т.е. для $z_1 < z_2$ всегда $\mu(z_1) \geq \mu(z_2)$.

К *строго монотонному типу* (SM -типу) отнесем методы непрерывного типа (C -типа) со скалярной функцией $Q(y) = f(y)$, в которых $\mu(z)$ — монотонно убывающая функция z , т.е. для $z_1 < z_2$ всегда $\mu(z_1) > \mu(z_2)$.

Приведенная классификация активно используется в теории T -представимых методов. Прежде чем перейти к ее изложению приведем примеры классификации методов, рассмотренных в разделах 4.3, 4.4.

Информационно–статистический алгоритм

Минимизируется липшицева функция f на отрезке $D = [a; b] \subset R^1$, т.е. $f \in \Phi = Lip[a; b]$. Если упорядочить нижним индексом последовательность координат выполненных испытаний Y_k ,

$$a = y_1 < y_2 < \dots < y_k = b,$$

то согласно (4.34) этому методу будет соответствовать функционал, вид которого для $y \in [y_{i-1}; y_i]$ определяется соотношением

$$W_k(\omega(f, Y_k), y) = \tilde{L}(y_i - y_{i-1}) \left(1 + (f_i - f_{i-1})^2 / (\tilde{L}(y_i - y_{i-1}))^2 \right) - 2(f_i + f_{i-1}) - 4(2y - (y_i - y_{i-1}) + (f_i - f_{i-1}) / \tilde{L})^2 / (y_i - y_{i-1}). \quad (4.80)$$

где \tilde{L} — оценка константы Липшица.

Нетрудно видеть, что определение 4.4 выполнено, а из (4.77) легко находится вид предельного функционала

$$W^*(y, f(y)) = -4f(y). \quad (4.81)$$

Таким образом, метод относится к строго монотонному типу (SM -типу) с функцией $\mu(z) = -4z$.

Метод С. А. Пиявского (многомерный вариант метода ломаных)

Минимизируется липшицева функция f в многомерном параллелепипеде $D \in R^N$, $N \geq 1$, т.е. $f \in \Phi = Lip(D)$ с константой L . Метод представим в форме (4.75) с функционалом

$$W_k(\omega(f, Y_k), y) = -\max \{ f^i - \tilde{L} \|y - y^i\| : y^i \in Y_k \} = \min \{ -f^i + \tilde{L} \|y - y^i\| : y^i \in Y_k \}, \quad (4.82)$$

где \tilde{L} — оценка константы Липшица.

Для него выполняется определение T -представимости с предельным функционалом

¹ Выбор обозначения соответствует принятому использованию символа C для именованного класса непрерывных функций.

$$W^*(y, f(y)) = -f(y). \quad (4.83)$$

Поэтому, по приведенной классификации, метод также относится к строго монотонному типу (SM – типу) с $\mu(z) = -z$.

Байесовский метод Х.Кушнера

Минимизируется непрерывная на $D=[a,b] \subset R^1$ функция f , т.е. $f \in \Phi = C[a,b]$. В качестве ее вероятностной модели используется винеровский случайный процесс, описанный в пункте 1.4.2.2 главы 1. Построенный на основе этой модели с использованием функции эффективности (4.35) одношагово–оптимальный метод, согласно (4.36), (1.57), (1.58), может быть представлен в нужной нам форме (4.75) с функционалом вида

$$W_k(\omega(f, Y_k), y) = (f_k^* - \delta_k - f_{i-1})/(y - y_{i-1}) + (f_k^* - \delta_k - f_i)/(y_i - y) \quad (4.84)$$

для $y \in [y_{i-1}, y_i]$, где невозрастающая последовательность $\delta_k \rightarrow \delta > 0$ при $k \rightarrow \infty$.

В качестве упражнения предлагается показать, что $\forall y \in D$ предел в (4.77) существует и, как это утверждалось в начале раздела 4.5, постоянен и равен $-\infty$, т.е.

$$W^*(y, f(y)) = const = -\infty.$$

Кроме того, выполняются требования Γ –представимости.

Таким образом, метод Х.Кушнера относится к константному типу.

Модифицированный метод Х.Кушнера для задач с известным значением f^*

Минимизируется непрерывная на $D=[a,b] \subset R^1$ функция f , гладкая в окрестности глобального минимума, т.е. $f \in \Phi = C[a,b] \cap C^1(O(Y^*))$. Будем дополнительно считать, что глобально–оптимальное значение функции f известно и равно f^* . Это значит, что существует индикатор $I(y)$ множества решений.

Модифицируем метод (4.75), (4.84) так, чтобы вместо оценки $f_k^* - \delta_k$ глобально–оптимального значения использовалось точное значение f^* . Новый метод при выборе очередной точки испытания в (4.75) будет использовать функционал вида

$$W_k(\omega(f, Y_k), y) = (f^* - f_{i-1})/(y - y_{i-1}) + (f^* - f_i)/(y_i - y). \quad (4.85)$$

Из гладкости $f(y)$ в окрестности Y^* следует, что значение $W^*(y, f(y)) = 0$ при $y \in Y^*$ и $W^*(y, f(y)) = -\infty$ при $y \notin Y^*$. Таким образом, согласно введенной классификации метод будет относиться к δ –типу.

Метод глобальной оптимизации для задач с ограничениями на основе адаптивных вероятностных моделей

В параллелепипеде $D \subset R^N$ ищется глобальный минимум функции f при наличии ограничений–неравенств $g_i(y) \leq 0$ ($i=1, \dots, m$).

Вектор-функция $Q=(f, g_1, \dots, g_m)$ должна принадлежать классу функций Φ'' со следующими свойствами. f, g_1, \dots, g_m — кусочно–непрерывны в D , множества $Y_{g_i} = \{y \in D : g_i(y) \leq 0\}$, ($i=1, \dots, m$) замкнуты, и при стремлении y изнутри Y_{g_i} к его границе функция $g_i(y)$ стремится к нулю, в любой окрестности решения $y^* \in Y^*$ существуют внутренние точки допустимого множества Y , и в пересечении Y с достаточно малой окрестностью y^* функция f непрерывна.

Под кусочной непрерывностью функции D в данном случае понимается ее непрерывность всюду в D за исключением конечного числа многообразий простой структуры размерности меньшей N , на которых допускаются конечные разрывы, равномерно ограниченные по величине.

Согласно (4.42), (1.62), (1.64), (1.68)–(1.70) функционал, определяющий выбор очередной точки испытания, имеет следующий вид.

$$W_k(\omega(f, g, Y_k), y) = F \left(\frac{f_k^* - \delta_k - M^f(\omega(f, Y_k), y)}{\sigma_k^f \min \{ \|y - y^i\|^s : y^i \in Y_k \}} \right) \times \\ \times F \left(\frac{-\varepsilon^G - \max \{ M^{g_i}(\omega(g_i, Y_k), y) : i = 1, \dots, m \}}{\sigma_k^G \min \{ \|y - y^i\|^\gamma : y^i \in Y_k \}} \right), \quad (4.86)$$

где $0 < F(z) < 1$ — функция распределения с нулевым средним и единичной дисперсией, f_k^* определяется формулой (4.40), $M^f(\cdot)$ — соотношением (1.63), (1.65), а $M^{g_i}(\cdot)$ — (1.67), $\gamma > s > 0$, $\sigma_k^f > 0$, $\sigma_k^G > 0$, $\delta_k \geq \delta > 0$, $\varepsilon^G > 0$.

При вычислении предела (4.77) для недопустимых точек $y \notin Y$ числитель в аргументе второго сомножителя в (4.86) всегда будет строго отрицателен, а знаменатель стремиться к нулю. С учетом того, что первый сомножитель ограничен единицей, получим $W^*(y, f(y), g(y)) = 0$. Если же точка y допустима, т.е. $y \in Y$, то при взятии предела (4.77) числитель первого сомножителя в (4.86) будет отрицателен, в силу правила выбора $f_k^* = f_k^*(\omega(f, g, Y_k))$ в (4.40) и условия на параметр $\delta_k \geq \delta > 0$, а знаменатель — стремиться к нулю. Следовательно, предельное значение всего выражения будет нулевым.

Итак, мы видим, что для предельного функционала выполнится

$$W^*(y, f(y), g(y)) = \text{const} = 0,$$

т.е. метод относится к группе методов константного типа (*Const* – типа).

4.5.3. Теория сходимости T–представимых алгоритмов

После рассмотрения примеров вернемся к изучению теории T–представимых методов.

Уточним понятие сходимости. Будем говорить, что *итерационный процесс* (4.75) *сходится на классе функций* Φ , если можно указать правило E_k , строящее по последовательности испытаний Y_k такую подпоследовательность

$$y^*(k) = E_k(\omega(Q, Y_k)) \in Y_k,$$

что все ее предельные точки являются решениями поставленной задачи, т.е. принадлежат Y^* .

Наряду с классификацией типов предельного функционала в теории T–представимости введем следующую типизацию характера сходимости последовательностей испытаний, согласующуюся с рассмотренными видами сходимости характеристически–представимых алгоритмов.

К *первому типу сходимости* отнесем всюду плотную сходимость (подчеркнем, что это не означает равномерного распределения точек испытаний в области поиска).

Для дальнейшего анализа всюду плотной сходимости удобно ввести специальный класс функций Φ' .

Описание класса Φ' . Этот класс включает такие вектор–функции Q , заданные в D и определяющие постановку решаемой задачи (пункт 4.5.1), для которых можно указать правило E_k , порождающее для любой всюду плотной в пределе на D последовательности (4.76) подпоследовательность $u^*(k) = E_k(\omega(Q, U_k)) \in U_k$, имеющую в качестве своих предельных точек только точки из множества решений Y^* .

Ко *второму типу сходимости* отнесем такое поведение, когда все предельные точки порожденной методом поисковой последовательности Y_∞ принадлежат множеству решений Y^* .

Кроме того, в задачах минимизации функций f на D без дополнительных функциональных ограничений выделим дополнительные свойства в поведении последовательности точек испытаний. Скажем, что ее поведение характеризуется *свойством приоритета*, если для двух любых предельных точек y' и y'' последовательности Y_∞ выполняется $f(y') = f(y'') \leq f(y^k)$ при всех $k \geq 1$. Скажем, что поведение характеризуется свойством *локальной сходимости* на D , если предельными точками множества Y_∞ могут являться только точки локальных минимумов функции $f(y)$ на D .

Основные результаты теории T –представимости связаны с выяснением влияния типа предельного функционала на поведение последовательности поисковых испытаний. Все утверждения относятся к классу задач, описанному в пункте 4.5.1.

Исследование сходимости первого и второго типов

Теорема 4.6. (необходимое условие всюду плотной сходимости).

Если T –представимый метод поиска решения задачи с неизвестным индикатором $I(y)$ множества решений Y^ обладает на классе функций Φ всюду плотной сходимостью, то его предельный функционал для $Q \in \Phi$ является константой, т.е. метод принадлежит $Const$ –типу на Φ .*

ДОКАЗАТЕЛЬСТВО. Пусть утверждение неверно, и $\exists Q \in \Phi \ni y', y'' \in D$, что

$$c' = W^*(y', Q(y')) < W^*(y'', Q(y'')) = c''.$$

Из всюду плотного типа сходимости следует, что y' и y'' — предельные точки для Y_∞ . Согласно (4.77) для произвольного c , удовлетворяющего неравенству $c' < c < c''$, $\exists K$, что $\forall k > K$

$$W_k(\omega(Q, Y_k), y'') > c.$$

С другой стороны, для y' существует сходящаяся к ней подпоследовательность испытаний y^{i_k} . По свойству T –представимости (4.78) получаем, что при $k \rightarrow \infty$

$$W_{i_k-1}(\omega(Q, Y_{i_k-1}), y^{i_k}) \rightarrow c' < c.$$

Но вместе с предыдущим неравенством это противоречит правилу (4.75) выбора новых точек испытаний. Теорема доказана.

Теперь выделим важный подкласс методов. Введем функцию расстояния до множества точек U_k

$$\rho(y, U_k) = \min \{ \|y - u\| : u \in U_k \}.$$

Определение 4.5. T – представимый метод назовем *Const* – нормальным, если для всякого U_k из (4.76) и $y \in D$ при $\rho(y, U_k) \geq R > 0$ выполняется условие отделимости от предельного значения

$$W_k(\omega(Q, U_k), y) \geq M > W^*(y, Q(y)), \quad (4.87)$$

где M может зависеть от Q, R и y .

Например, для метода Х.Кушнера из (4.84) непосредственно видно, что числители двух дробей, входящих в функционал W_k всегда больше или равны $f^* - \delta_0 - f^{\max} < 0$, где f^{\max} — верхняя оценка f на $[a, b]$, существующая в силу непрерывности f , а знаменатели при $\rho(y, U_k) \geq R$ больше $R > 0$. Отсюда вытекает, что

$$W_k(\omega(f, U_k), y) \geq M = -(f^{\max} - f^* + \delta_0) / R > W^* = -\infty.$$

Таким образом, метод Х.Кушнера является K – нормальным.

Контрольные вопросы и упражнения

Покажите, что при определенных дополнительных требованиях на класс функций $Q=(f, g)$ метод условной глобальной оптимизации на основе адаптивных стохастических моделей, определяемый функционалом (4.86), будет *Const* – нормальным. Укажите эти требования.

Теорема 4.7. (критерий всюду плотной сходимости). На функциях класса $\Phi \subseteq \Phi'$ в задачах с неизвестным индикатором множества решений, описанных в пункте 4.5.1, *Const* – нормальный метод поиска обладает всюду плотным типом сходимости тогда и только тогда, когда он принадлежит *Const* – типу, т.е. $W^*(y, Q(y)) = const$.

НЕОБХОДИМОСТЬ непосредственно вытекает из теоремы 4.6.

ДОСТАТОЧНОСТЬ. Пусть $W^*(y, Q(y)) \equiv c$, но $\exists y' \in D$, не являющаяся предельной для последовательности испытаний Y_∞ . Тогда найдется $R > 0$, что для всех k $\rho(y, Y_k) \geq R$. Из (4.87) следует, что $\exists M$, что $\forall k$

$$W_k(\omega(Q, Y_k), y') \geq M > c.$$

Но т.к. D — компакт, в D существует по крайней мере одна предельная точка y'' для Y_∞ . На сходящейся к ней подпоследовательности y^{i_k} при $k \rightarrow \infty$

$$W_{i_k-1}(\omega(Q, Y_{i_k-1}), y^{i_k}) \rightarrow c.$$

Но это противоречит правилу выбора точек испытаний (4.75) с учетом предыдущего неравенства. Таким образом, метод размещает испытания всюду плотно, а т.к., по предположению, класс Φ включается в описанный в начале пункта 4.5.3 класс Φ' , то метод будет обладать на нем всюду плотной сходимостью. Теорема доказана.

Заметим, что из теоремы 4.7 и рассмотренного выше следует, что на соответствующих им классах функций метод Х.Кушнера и метод условной глобальной оптимизации на основе адаптивных стохастических моделей сходятся со всюду плотным характером сходимости. Как уже отмечалось ранее, всюду плотный характер сходимости не означает равномерного размещения точек испытаний в области поиска. Их концентрация в усеченной последовательности Y_k , выражаемая расстояниями между ближайшими точками испытаний, может

быть существенно различной в различных подобластях области поиска и даже иметь различный порядок малости при $k \rightarrow \infty$. Соответствующие оценки относительной концентрации будут получены в разделе 4.6.

Полученные выше условия сходимости выполняются при весьма слабых предположениях о решаемой задаче, когда невозможно построить нижние оценки значений функций. Соответственно и сходимость достигается только за счет всюду плотного в пределе размещения испытаний. Чтобы изменить тип сходимости, нужна дополнительная информация о задаче. Рассмотрим случай, когда известен индикатор $I(y)$ множества решений Y^* (при этом, конечно, до проведения поиска нельзя указать точку $y^* \in Y^*$). Индикатор $I(y)$ является известным, например, в тех задачах глобальной оптимизации, к которым сводятся задачи поиска решений систем совместных уравнений и неравенств. Таким задачам адекватны T – представимые методы δ – типа.

Введем понятие δ – нормальности, аналогичное по смыслу $Const$ – нормальности из определения 4.5, но ориентированное на применение к методам δ – типа.

Определение 4.6. T – представимую процедуру поиска назовем δ – нормальной, если для $y \in D$ и всякого множества U_k из (4.76) при $\rho(y, U_k) \geq R > 0$ выполняется

$$W_k(\omega(Q, U_k), y) \geq M > \min\{W^*(y, Q(y)) : y \in D\}, \quad (4.88)$$

где M может зависеть от Q, R и y , но не зависит от k .

Изучим условия сходимости второго типа.

Теорема 4.8. (достаточные условия сходимости второго типа). Если на классе функций Φ метод поиска является δ – нормальным и относится к δ – типу, то он порождает для задач с непустым замкнутым множеством решений Y^* последовательности поисковых испытаний, все предельные точки которых принадлежат Y^* . Если же $Y^* = \emptyset$, метод порождает всюду плотную в пределе последовательность испытаний.

Доказательство. Если $Y^* = \emptyset$, то δ – нормальный метод δ – типа является в гиперпараллелепипеде D $Const$ – нормальным методом константного типа. Поэтому, из доказательства теоремы 4.7 вытекает его всюду плотная сходимость.

Пусть теперь Y^* замкнуто и непусто. Если бы никакая точка y' из Y^* не была предельной для последовательности испытаний, то рассуждая аналогично тому, как это было сделано в доказательстве теоремы 4.7, получили бы, что в области $D \setminus Y^*$ точки испытаний располагались в пределе всюду плотно. Но в этом случае найдутся граничные точки Y^* , являющиеся предельными для последовательности испытаний. Однако в силу замкнутости Y^* это противоречит первоначальному предположению, следовательно, в Y^* существует предельная точка y' .

Из (4.77) и определения сходимости δ – типа при $k \rightarrow \infty$ имеем

$$W_k(\omega(Q, Y_k), y') \rightarrow \eta(1) > \eta(0). \quad (4.89)$$

Следовательно, ни одна из точек $y'' \in D \setminus Y^*$ не может быть предельной, т.к. иначе на сходящейся к ней подпоследовательности y^{i_k} при $k \rightarrow \infty$ из (4.78) имели бы

$$W_{i_{k-1}}(\omega(Q, Y_{i_{k-1}}), y^{i_k}) \rightarrow \eta(0),$$

но это вместе с (4.89) противоречит правилу выбора точек испытаний (4.75). Теорема доказана.

Исследование сходимости со свойством приоритета и локальной сходимости

Перейдем теперь к изучению введенных в начале пункта 4.5.3 свойств приоритета и локальной сходимости применительно к более частной задаче поиска глобального минимума на компакте D без дополнительных функциональных ограничений

$$f(y) \rightarrow \min, y \in D \subset R^N. \quad (4.90)$$

Впервые эти свойства были выявлены и исследованы Р.Г.Стронгиным [7] для конкретных информационно–статистических алгоритмов глобального поиска, а затем в работах В.А.Гришагина теоретически обобщены для класса характеристических алгоритмов оптимизации (см. раздел 4.4).

Более общие результаты в предположении непрерывности целевой функции f в задаче (4.90) формулирует следующая теорема.

Теорема 4.9. (необходимые и достаточные условия свойства приоритета)
 Для того, чтобы T –представимый метод поиска минимума в задаче (4.90) для скалярных непрерывных в компакте D функций f породил поисковую последовательность со свойством приоритета, достаточно принадлежности метода строго монотонному типу (SM –типу), а на классе методов непрерывного типа (C –типа) — необходимо, чтобы метод поиска относился к монотонному типу (M –типу).

ДОСТАТОЧНОСТЬ. Пусть y' и y'' — предельные точки последовательности испытаний Y_∞ и $f(y') \neq f(y'')$, например, $f(y') < f(y'')$. Из строго монотонного типа метода следует, что $\forall y \in D W^*(y, f(y)) = \mu(f(y))$ и $\mu(f(y')) > \mu(f(y''))$. Из (4.77), (4.79) и непрерывности $\mu(z)$ и $f(y)$ сразу получаем, что при достаточно малом $\varepsilon > 0$ найдется точка испытания y^{k*} , достаточно близкая к y' , что $\forall k > k_*$

$$W_k(\omega(f, Y_k), y^{k*}) > W^*(y', f(y')) - \varepsilon > W^*(y'', f(y'')),$$

но тогда y'' не может быть предельной для множества точек испытаний. Полученное противоречие доказывает, что $f(y') = f(y'')$.

Аналогично получаем, что $f(y') \leq f(y^k) \quad \forall y^k \in Y_\infty$. Таким образом для поисковой последовательности выполняется свойство приоритета.

НЕОБХОДИМОСТЬ. Предположим, что метод относится к C –типу, но не принадлежит монотонному типу. Тогда найдутся $z_1 < z_2$, что $\mu(z_1) < \mu(z_2)$. Построим непрерывную функцию $f(y)$, чтобы $f(y^1) = z_1$ и $\forall y \in D f(y) \geq z_1$. Пусть метод по начальной информации $\omega_1 = \omega(f, \{y^1\})$ строит следующую точку измерения y^2 . Всегда можно выбрать функцию $f(y)$ так, чтобы $f(y^2) = z_2$.


В силу компактности области D у последовательности точек испытаний Y_∞ существует хотя бы одна предельная точка y' . Рассмотрим значение $f(y')$. Из-за свойства приоритета $f(y') \leq f(y^1)$. Но, по построению функции, $f(y) \geq f(y^1) \quad \forall y \in D$, следовательно $W^*(y', f(y')) = \mu(f(y')) = \mu(f(y^1)) < \mu(f(y^2))$. Отсюда, в силу правила (4.75), получаем противоречие. Следовательно, никакая точка y' не может быть предельной для Y_∞ . Поскольку это невозможно, метод должен быть монотонен. Таким образом, доказательство завершено.

Заметим, что из условий теоремы 4.9 и ранее проведенного анализа следует, что информационно–статистический метод и метод С.А. Пиявского обладают свойствами приоритета на классе липшицевых функций, а метод Х.Кушнера этим свойством не обладает.

Следует обратить внимание на полезное следствие, вытекающее из свойства приоритета.

Следствие 4.9.1. Если T – представимый метод поиска решения задачи (4.90) на компакте D сходится, то для непрерывных в D функций f из выполнения свойства приоритета следует конечность числа испытаний, размещаемых методом в любой подобласти со значениями $f(y) \geq f(y^*) + \varepsilon$ при $\varepsilon > 0$.

ДОКАЗАТЕЛЬСТВО. В самом деле, если процедура сходится (в смысле понятия сходимости, введенного в начале пункта 4.5.3), то в Y^* имеются предельные точки последовательности Y_∞ . Отсюда и из свойства приоритета очевидно следует отсутствие предельных точек последовательности Y_∞ в замкнутых подобластях, не содержащих точек из Y^* , т.е. конечность числа точек испытаний в таких подобластях. Доказательство завершено.

 **Замечание.** Гарантировать сходимость к глобальному минимуму в сочетании со свойством приоритета можно только на классах функций, допускающих конечные оценки снизу для f по конечному числу испытаний.

ДОКАЗАТЕЛЬСТВО. Если конечных минорант для f не существует, то глобальный минимум может располагаться в любой точке, не совпадающей с точками проведенных испытаний. Тогда для его определения необходимо всюду плотное в пределе размещение испытаний. Но это противоречит следствию 4.9.1. Доказательство завершено.

В дальнейшем предполагаем, что $f \in \Phi_q$, где класс Φ_q включает непрерывные функции, которые для любого U_k из (4.76) допускают конечную нижнюю оценку вида

$$f(y) \geq q(\omega(f, U_k), y). \quad (4.91)$$

Если данная оценка хотя и существует, но точно неизвестна, гарантировать для процедуры поиска сходимость к множеству глобальных минимумов при выполнении свойства приоритета заранее нельзя, и условие сходимости необходимо носит апостериорный характер. В ходе поиска оно может и не выполниться. Процедура поиска должна сохранять при этом по крайней мере локальную поисковую способность, которая отражается свойством локальной сходимости. Выполнение этого свойства требует от характеристического функционала $W_k(\cdot)$ некоторых специальных свойств. Пусть

$$D_y = \{u \in D : f(u) < f(y)\}. \quad (4.92)$$

Рассматривая вопрос на содержательном уровне, можно сказать, что функционал $W_k(\cdot)$ должен позволять идентифицировать наличие не содержащих точек испытаний подобластей D_y , когда точка y , принадлежащая границе D_y , является точкой сгущения множества точек испытаний. Поскольку вид областей D_y может быть весьма разнообразным, удобно ввести некоторый "стандартный"

набор Γ пар $(A; y)$, где A — множество, $y \in \partial A$ (∂A — граница A) и рассматривать указанное выше свойство, заменяя D_y на $A \subseteq D_y$.

Определение 4.7. Пусть Γ — некоторая совокупность пар $(A; y)$, где $y \in D$, а A — непустые, открытые, ограниченные, связные подмножества из D , причем $y \in \partial A$. T — представимый метод решения задачи (4.90) назовем M — нормальным (относительно класса функций Φ и класса Γ пар $(A; y)$), если для любых $f \in \Phi$ и $y \in D$ таких, что существует пара $(A; y) \in \Gamma$ с $A \subset D_y$ из (4.92), для произвольной последовательности U_∞ , имеющей y в числе своих предельных точек и такой, что $U_\infty \cap D_y = \emptyset$, начиная с некоторого k выполняется

$$\sup\{W_k(\omega(f, U_k), u) : u \in A\} \geq M > W^*(y, f(y)). \quad (4.93)$$

При этом M может зависеть от f и $(A; y)$, но не зависит от k .

Укажем условия, при которых можно рассчитывать на наличие у поисковой последовательности свойств локальной сходимости и приоритета. Сформулируем их в виде следующей теоремы.

Теорема 4.10. (условия локальной сходимости со свойством приоритета) Пусть $f \in \Phi \subseteq \Phi_q$ из (4.91), число локальных минимумов f в D конечно, метод поиска решения задачи (4.90) относится к строго монотонному типу (SM -типу) с $W^*(y, f(y)) = \mu(f(y))$ и является M — нормальным относительно Φ и класса пар Γ . Тогда справедливы следующие утверждения.

A. Если y' — предельная точка последовательности испытаний Y_∞ , то для любого $k \geq 1$ $f(y') \leq f(y^k)$ и либо точка y' является локальным минимумом f на D , либо для y' в классе Γ не существует пар $(A; y') \in \Gamma$ со сколь угодно малыми значениями диаметра $\text{diam}(A)$ и $A \subset D_{y'}$.

B. Если y'' — любая другая предельная точка последовательности испытаний Y_∞ , то $f(y'') = f(y')$.

Доказательство. Справедливость утверждения *B* и первой части утверждения *A* непосредственно следует из теоремы 4.9. Доказательство остальных утверждений читатель сможет найти в работе [15].

Анализ характера сходимости к локально-оптимальным точкам

Приведенная теорема 4.10, хотя и характеризует особенности возможных предельных точек последовательности испытаний, но не дает ответа на вопрос о том, как именно происходит приближение точек сходящихся подпоследовательностей к их предельным значениям.

Теория T — представимых методов позволяет проанализировать характер этого приближения. Для возможности такого анализа рассмотрим семейства множеств $K_\varepsilon(v, \varphi)$, являющихся пересечениями ε — окрестностей $O_\varepsilon(0)$ точки 0 с конусами $K_\varepsilon(v, \varphi)$ следующего вида

$$K(v, \varphi) = \{y \in R^N : (v, y) > \|y\| \cos \varphi\},$$

где $\|v\| = 1$.

Введем также класс $\Psi(y, K_\varepsilon(v, \varphi))$ последовательностей U_k попарно различных точек испытаний из (4.76), что y будет предельной для U_∞ и $K_\varepsilon(v, \varphi) \cap U_\infty = \emptyset$.

Определение 4.8. Пусть выбрано фиксированное $\varphi^* \in (0, \pi/2]$. T – представимый метод решения задачи (4.90) с функцией f из класса Φ назовем φ – нормальным, со значением $\varphi = \varphi^*$ на классе Φ , если для всякой внутренней по отношению к D локально–оптимальной точки y^0 задачи (4.90) $\forall v$ с $\|v\|=1$, $\forall \varepsilon > 0$ и произвольной последовательности точек испытаний $U_\infty \in \Psi(y^0, K_\varepsilon(v, \varphi))$, удовлетворяющей дополнительному требованию, что $\forall u_k \in U_\infty: f(u_k) \geq f(y^0)$, начиная с некоторого шага k (возможно, зависящего от U_∞) выполняется следующее неравенство

$$\sup \{ W_k(\omega(f, U_k), u): u \in y^0 + K_\varepsilon(v, \varphi) \} \geq M(f, y) > W^*(y, f(y)). \quad (4.94)$$

Необходимые результаты представлены в следующей теореме.

Теорема 4.11 (условия «сходимости» по секторам). Пусть T – представимый метод решения задачи (4.90) с функцией f из класса Φ относится к строго монотонному типу (SM –типу) и при некотором $\varphi = \varphi^* \in (0, \pi/2]$ является φ – нормальным на классе Φ . Тогда, если некоторая внутренняя, по отношению к области D , локально–оптимальная точка y^0 будет предельной для последовательности испытаний, то для любого направления v ($\|v\|=1$) в коническом секторе $y^0 + K(v, \varphi^*)$ с угловым размером $2\varphi^*$ и ориентацией v найдется подпоследовательность испытаний, сходящаяся к y^0 .

Доказательство. Пусть внутренняя для D локально–оптимальная точка y^0 является предельной для последовательности испытаний Y_∞ , но утверждение теоремы для нее не верно. Тогда найдется конический сектор $y^0 + K(v, \varphi^*)$ с угловым размером $2\varphi^*$ и ориентацией v ($\|v\|=1$), что $(y^0 + K(v, \varphi^*)) \cap Y_\infty = \emptyset$. Поскольку метод строго монотонен, то по теореме 4.9 для него выполнено свойство приоритета, т.е., в частности, для любой точки испытания y^k выполнено $f(y^0) \leq f(y^k)$ ($k=1, 2, \dots$).

Таким образом, для последовательности испытаний Y_∞ выполнены все условия, накладываемые на множества U_∞ в определении 4.8 и, следовательно, $\forall \varepsilon > 0$ последовательность Y_∞ принадлежит классу $\Psi(y^0, K_\varepsilon(v, \varphi^*))$.

Отсюда и из свойства φ^* – нормальности следует, что найдутся такое значение $M(f, y^0)$ и точка $y_\varepsilon \in (y^0 + K_\varepsilon(v, \varphi^*))$ что, начиная с некоторого шага $k=k(\varepsilon)$, будут выполняться неравенства

$$W_k(\omega(f, Y_k), y_\varepsilon) \geq M(f, y^0) > W^*(y^0, f(y^0)) = \mu(f(y^0))$$

С другой стороны, на сходящейся к y^0 подпоследовательности y_{k_s} при $s \rightarrow \infty$, в силу (4.78),

$$W_{k_s-1}(\omega(f, Y_{k_s-1}), y_{k_s}) \rightarrow W^*(y^0, f(y^0)) = \mu(f(y^0)),$$

и $\forall \varepsilon: 0 < \varepsilon < 0.5 \cdot (M(f, y^0) - \mu(f(y^0)))$ при некотором достаточно большом k_s и дополнительном условии $k_s > k(\varepsilon)$ окажется, что

$$W_{k_s-1}(\omega(f, Y_{k_s-1}), y^{k_s}) < W^*(y^o, f(y^o)) + \varepsilon < W_{k_s}(\omega(f, Y_{k_s}), y_\varepsilon),$$

но это противоречит правилу выбора точек испытаний (4.75). Таким образом, теорема верна.

Исследование сходимости только к глобально-оптимальным точкам

Наибольший интерес представляют условия, при которых для T –представимых процедур строго монотонного типа (SM –типа) все предельные точки последовательности испытаний являлись бы глобально–оптимальными. Как уже отмечалось выше, эти условия неизбежно включают в себя выражение для миноранты целевой функции.

Теорема 4.12 (условия глобальной оптимальности предельных точек). Пусть $f \in \Phi \subseteq \Phi_q$ из (4.91), число локальных минимумов f в D конечно, T –представимый метод поиска решения задачи (4.90) относится к строго монотонному типу (SM –типу) с $W^*(y, f(y)) = \mu(f(y))$. Тогда справедливы следующие утверждения.

A. Если начиная с некоторого k выполняются неравенства

$$W_k(\omega(f, Y_k), y) \geq \mu(q(\omega(f, Y_k), y)), \quad (4.95)$$

то все предельные точки поисковой последовательности принадлежат множеству глобальных минимумов Y^* .

B. Если, начиная с некоторого k , для любого $y \in D$ при $\rho(y, Y_k) \geq R > 0$ существует M , что выполняется оценка

$$W_k(\omega(f, Y_k), y) \geq M(R, f, y) > \mu(q(\omega(f, Y_k), y)), \quad (4.96)$$

где $M(\cdot)$ не зависит от k но может зависеть от R, f и y , то предельными точками последовательности испытаний являются все точки множества глобальных минимумов Y^* и только они.

ДОКАЗАТЕЛЬСТВО приведено в [15].

Контрольные вопросы и упражнения.

Самостоятельно докажите утверждения *A* и *B* теоремы 4.12.

Пример. Проиллюстрируем применение теорем 4.10–4.12 к анализу свойств метода С.А. Пиявского. Рассмотрим его на классе липшицевых функций для размерности $N > 1$. Этот случай интересен тем, что его нельзя проанализировать с позиций теории характеристически представимых алгоритмов, поскольку область ее применимости ограничивается задачами с одним переменным и сводимыми к ним.

Ранее было показано, что при $N \geq 1$ метод С.А. Пиявского, определяемый соотношением (4.82), является T –представимым и относится к строго монотонному типу (SM –типу) с $\mu(z) = -z$. Предположим, что в методе (4.82) оценка константы Липшица \tilde{L} выбрана так, что $\tilde{L} \geq L$.

ПРОВЕРКА УСЛОВИЯ M –НОРМАЛЬНОСТИ. В качестве класса Γ выберем всевозможные пары $(A; y)$, где A — произвольные открытые односвязные множества, а $y \in \partial A$. Рассмотрим точку, не являющуюся локальным минимумом f на D , тогда D_y из (4.92) не пусто и существует $(A; y) \in \Gamma$, что $A \subseteq D_y$.

Рассмотрим произвольную точку $\bar{u} \in A$. Поскольку она для A внутренняя, то в условиях определения 4.7 $\exists R > 0$, что $\forall k \rho(\bar{u}, U_k) \geq R > 0$. Из условия Липшица следует, что $f(u^i) \leq f(\bar{u}) + L\|\bar{u} - u^i\|$. Используя это неравенство совместно с (4.82) и учитывая, что $f(\bar{u}) < f(y)$, получим оценку

$$\begin{aligned} & \sup\{W_k(\omega(f, U_k), u): u \in A\} \geq W_k(\omega(f, U_k), \bar{u}) = \\ & = \min\{-f(u^i) + \tilde{L}\|\bar{u} - u^i\|: i = 1, \dots, k\} \geq \min\{-f(\bar{u}) + (\tilde{L} - L)\|\bar{u} - u^i\|: i = 1, \dots, k\} \geq \\ & \geq -f(\bar{u}) + (\tilde{L} - L)R > -f(y) = W^*(y, f(y)). \end{aligned}$$

Таким образом, доказана M -нормальность по отношению к введенному классу пар Γ .

Заметим, что рассматриваемый класс пар содержит пары $(A; y)$ с множествами A сколь угодно малого диаметра. Поэтому из строгой монотонности метода, доказанной M -нормальности и теоремы 4.10 следует, что процесс поиска в методе С.А. Пиявского при выполнении условия $\tilde{L} \geq L$ обладает свойством приоритета и локальной сходимости.

АНАЛИЗ φ -НОРМАЛЬНОСТИ ПРИ $\tilde{L} > L$. В указанных условиях проверим возможность выполнения для метода С.А. Пиявского оценки (4.94). Рассмотрим произвольную точку y^0 внутреннего локального минимума f на D и некоторое направление v . Вначале выберем последовательности U_∞ из класса $\Psi(y, K_\varepsilon(v, \varphi))$, включающие только такие точки измерений u_k , которые располагаются на прямой $y^0 + vt$, $t \in \mathbb{R}^1$ вне промежутка $(0; \varepsilon)$ по t . При достаточно больших номерах шагов для них, с учетом существования в U_∞ подпоследовательности точек, сходящейся к y^0 , возможна (при достаточно малых $\delta > 0$) оценка

$$\begin{aligned} & \sup\{W_k(\omega(f, U_k), u): u \in K_\varepsilon(v, \varphi)\} \geq \\ & \geq -\delta + \max\left\{\min\{-f(y^0) + \tilde{L} \cdot t; -f(y^0) - L \cdot \varepsilon + \tilde{L} \cdot (\varepsilon - t)\}: 0 \leq t \leq \varepsilon\right\} = -f(y^0) + \tilde{L}\Delta. \end{aligned}$$

где $\Delta = \varepsilon(\alpha - 1)/(2\alpha)$, $\alpha = \tilde{L}/L$. Таким образом, для данного подкласса последовательностей U_∞ в качестве «худшей» в значении $\sup\{\dots\}$ можно принять точку $u = \bar{u} = y^0 + v \cdot \Delta$.

Нетрудно понять, что нижняя оценка для $W_k(\omega(f, U_k), \bar{u})$ не понизится, если в U_∞ добавить произвольные точки $u_k \in (y^0 + K(v, \varphi))$, но $u_k \notin O_\varepsilon(y^0)$.

Отсюда следует, что справедливость (4.94) можно теперь установить, оценивая снизу в выбранной точке \bar{u} точную нижнюю грань значения $W_k(\omega(f, U_k), \bar{u})$ по всем последовательностям U_∞ , имеющим предельную точку в y^0 и не пересекающимся с открытым коническим множеством $y^0 + K(v, \varphi)$.

Нетрудно видеть, что проверка выполнения (4.94) в указанных условиях эквивалентна доказательству того, что при соответствующем выборе $\beta = tg(\varphi)$ значение минимума в задаче

$$a^* = \min\left\{\alpha\sqrt{r^2 + (t - \Delta)^2} - \sqrt{r^2 + t^2} : r \geq \beta \cdot t \geq 0\right\}$$

строго положительно ($a^* > 0$). Здесь, как и ранее, $\alpha = \tilde{L}/L$, а $\Delta = \varepsilon(\alpha - 1)/(2\alpha)$. Переменные t и r являются проекциями вектора $y - y^0$ на направление v и ортогональное ему подпространство.

Казалось бы, необходимо дополнительно проверить выполнение такой же оценки $a^* > 0$ при взятии минимума по другой области, когда $r \geq 0, t \leq 0$, но этого можно не делать, т.к. данная оценка, очевидно, выполняется.

Заменой переменных $t := t, \rho = \sqrt{r^2 + t^2}$ интересующая нас задача приводится к виду

$$0 < a^* = \min \left\{ \alpha \sqrt{\rho^2 - 2 \cdot t \cdot \Delta + \Delta^2} - \rho : t \leq \gamma \cdot \rho \right\},$$

где $\gamma = 1/\sqrt{1 + \beta^2}$. Необходимо выяснить, при каких $\varphi = \varphi^*$ это неравенство справедливо. Для его выполнения изолиния с нулевым значением минимизируемой функции, удовлетворяющая уравнению

$$\rho^2(1 - 1/\alpha^2) - 2 \cdot t \cdot \Delta + \Delta^2 = 0,$$

не должна иметь пересечений с границей допустимой области $t = \rho \cdot \gamma$. Таким образом, задача сводится к обеспечению отрицательности дискриминанта

$$d = 4 \cdot \Delta^2 \gamma^2 - 4 \cdot (1 - 1/\alpha^2) \Delta^2 < 0.$$

Это неравенство эквивалентно следующему

$$\operatorname{tg} \varphi^* > 1/\sqrt{\alpha^2 - 1}.$$

Таким образом, метод С.А. Пиявского φ -нормален с $\varphi = \varphi^*$, где φ^* удовлетворяет приведенному выше неравенству.

Таким образом, например, при $\tilde{L} > \sqrt{2}L$ (т.е. при $\alpha = \sqrt{2}$) можно гарантировать, что при сходимости подпоследовательности испытаний к некоторому локальному минимуму y^o в любом конусе с угловым размером $2\varphi^* > \pi/2$ и вершиной в точке y^o найдется сходящаяся к y^o подпоследовательность измерений, целиком лежащая в этом конусе. Видно, что при $\alpha \rightarrow \infty$ угловой размер $2\varphi^*$ соответствующих конусов стремиться к нулю.

ПРОВЕРКА УСЛОВИЙ ГЛОБАЛЬНОЙ СХОДИМОСТИ. Очевидно, что при $\tilde{L} \geq L$ для метода С.А. Пиявского выполнено условие (4.95), для которого в (4.91) можно принять

$$q(\omega(f, Y_k), y) = \max \left\{ f(y^i) - L \|y - y^i\| : y^i \in Y_k \right\}.$$

В силу теоремы 4.12 отсюда следует, что при $\tilde{L} \geq L$ все предельные точки последовательности испытаний принадлежат множеству глобальных минимумов Y^* .

Если с некоторого шага оценка константы Липшица $\tilde{L} > L$, то выполнено усиленное условие (4.96), где при $\rho(y, Y_k) \geq R > 0$

$$M(R, f, y) = \min \left\{ -f(y^i) + \tilde{L} \|y - y^i\| : y^i \in Y_k \right\} \geq -q(\omega(f, Y_k), y) + (\tilde{L} - L)R > \\ > \mu(q(\omega(f, Y_k), y)).$$

Таким образом, при $\tilde{L} > L$ множество предельных точек последовательности испытаний метода С.А. Пиявского совпадает с Y^* . Рассмотрение примера закончено.

Итак, мы познакомились с теорией сходимости T -представимых методов и проиллюстрировали ее положения на примере нескольких методов многоэкстремальной оптимизации.

Относительно методов со всюду плотным характером сходимости остался не изученным вопрос о методах аналитического анализа относительной плотности размещения их испытаний. Эта тема обзорно рассматривается в следующем разделе.

4.6. Анализ относительной плотности размещения испытаний при всюду плотной сходимости

Аналитические методы анализа относительной плотности размещения испытаний являются необходимым инструментом исследования методов со всюду плотным характером сходимости. Сам факт всюду плотной сходимости ничем не отличает обладающие такой сходимостью методы от тривиального равномерного перебора, однако опубликованные результаты экспериментального исследования таких методов показывают их высокую эффективность [28, 33, 51]. Теория аналитического вывода оценок относительной плотности предложена в [15]. В этом разделе излагаются основы подхода на примере метода Х.Кушнера, а затем в обзорном плане — общая методика и результаты ее применения для ряда алгоритмов.

4.6.1. Необходимые понятия и обозначения

Воспользуемся еще раз определенным в (4.76) усечением U_k последовательности U_∞ произвольных попарно различных точек испытаний из D . Для $y' \in U_k$ введем для $N > 1$ функцию $\bar{\rho}(y', U_k)$ расстояния от y' до ближайшей точки из U_k не совпадающей с y' . В случае же размерности $N=1$ под $\bar{\rho}(y', U_k)$ будем понимать расстояние от $y' \in U_k$ до ближайшей справа точки из U_k . Если множество испытаний U_k порождено методом оптимизации типа (4.75), т.е. $U_k = Y_k$, будем использовать более короткую запись

$$\bar{\rho}_k(y') = \bar{\rho}(y', Y_k), \quad (4.97)$$

которая ниже будет часто применяться.

Эти обозначения иллюстрирует рис. 4.5, где жирными точками для $N=2$ и вертикальными штрихами для $N=1$ отмечены элементы множества Y_k . Показано различие в трактовке обозначения $\bar{\rho}(y)$ для разных размерностей.

Для всюду плотно сходящихся процедур величины $\bar{\rho}_k(y)$, $y \in Y_k$ являются бесконечно малыми при $k \rightarrow \infty$.

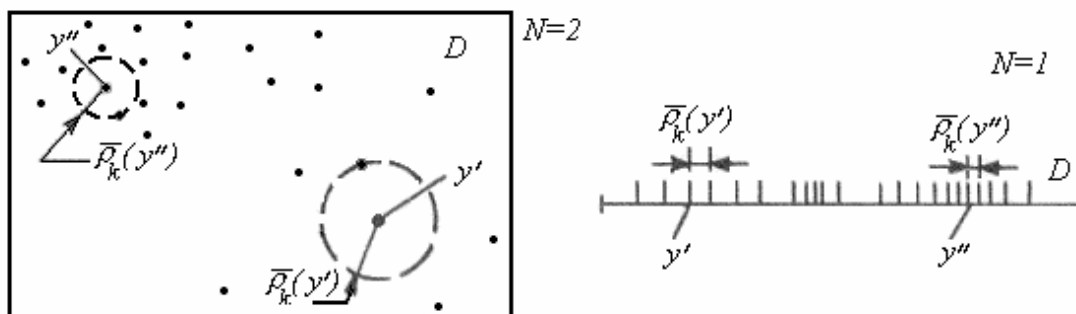


Рис. 4.5. Значения $\bar{\rho}_k(y')$ и $\bar{\rho}_k(y'')$ для размерности $N=2$ и $N=1$

Определение 4.8. Пусть y' и y'' — точки испытаний, тогда если при некотором $p > 0$ отношение

$$\alpha_k(y', y'') = (\bar{\rho}_k(y'')) / (\bar{\rho}_k(y'))^{1/p} \quad (4.98)$$

при $k \rightarrow \infty$ ограничено и строго положительно, то число p будем называть порядком относительной концентрации испытаний, а функцию $\alpha_k(y', y'')$ — относительной концентрацией (плотностью) испытаний в точке y' по отношению к точке y'' .

Для всюду плотно сходящихся методов при бесконечно долго продолжающемся процессе вычислений в любой окрестности всякой точки ранее выполненного испытания когда-либо обязательно проводится новое, поэтому $\forall y', y'' \in Y_k$ с ростом k будут встречаться моменты, когда станет уменьшаться величина $\bar{\rho}_k(y')$, а также — моменты, когда уменьшаться будет $\bar{\rho}_k(y'')$. Отсюда следует, что в процессе поиска относительная концентрация испытаний $\alpha_k(y', y'')$ при $k \rightarrow \infty$ будет совершать колебания в некоторых границах. Наша цель будет состоять в том, чтобы найти зависимость границ этих колебаний для произвольно выбранных точек y' и y'' из D от поведения функции Q , описывающей исходную задачу.

Пример упрощенного способа оценивания относительной концентрации

Прежде чем перейти к формальному описанию метода оценивания границ колебаний относительной плотности размещения испытаний, приведем простой способ получения некоторого «среднего» значения $\tilde{\alpha}(y', y'')$, вокруг которого совершает колебания относительная плотность при $k \rightarrow \infty$. Этот способ не является строгим, но позволяет лучше понять основную идею последующего формального изложения.

Покажем, как можно получить значение $\tilde{\alpha}(y', y'')$ на примере метода Х. Кушнера. Предположим, что порядок относительной концентрации испытаний p в (4.90) равен единице ($p=1$). Тогда

$$\bar{\rho}_k(y'') = \alpha_k(y', y'') \bar{\rho}_k(y').$$

Для метода Х.Кушнера с функционалом (4.36), как было указано в пункте 4.3.4, для отрезка $[y_{i-1}; y_i]$ можно вычислить характеристику

$$R(i) = \max \{W_k(\omega(f, Y_k), y) : y \in [y_{i-1}; y_i]\} = -4 \frac{(f_k^* - \delta_k - f(y_i))(f_k^* - \delta_k - f(y_{i-1}))}{y_i - y_{i-1}}.$$

Выберем и зафиксируем произвольную точку испытания y' . Пусть на k -м шаге она совпадает с $(i-1)$ -й точкой в порядке возрастания координат y_{i-1} , где $i=i(k, y')$. Тогда в наших обозначениях

$$y_i - y_{i-1} = \bar{\rho}_k(y').$$

Аналогично можно рассмотреть другую точку испытания y'' , причем для $j=i(k, y'')$

$$y_j - y_{j-1} = \bar{\rho}_k(y'') = \alpha_k(y', y'') \cdot \bar{\rho}_k(y').$$

В процессе поиска, из-за всюду плотного характера сходимости, значения $R(i(k, y'))$ и $R(i(k, y''))$ будут поочередно становиться то больше, то меньше одно

другого, поэтому «среднее» значение $\tilde{\alpha}(y', y'')$ можно найти, приравняв их друг к другу

$$R(i(k, y')) = R(i(k, y''))$$

с заменой в них значения $\alpha_k(y', y'')$ на «среднее» $\tilde{\alpha}(y', y'')$. Отсюда, используя непрерывность функции f и всюду плотный характер сходимости, при $k \rightarrow \infty$ получим

$$(f^* - \delta - f(y'))^2 = (f^* - \delta - f(y''))^2 / \tilde{\alpha}(y', y'')$$

Выберем теперь $y' = y$, $y'' = y^*$. Тогда

$$\tilde{\alpha}(y, y^*) = \left(\frac{1}{1 + (f(y) - f^*)/\delta} \right)^2 \quad (4.99)$$

Это выражение характеризует «среднее» значение вокруг которого колеблется плотность испытаний в точке y по отношению к плотности в точке глобального минимума y^* .

Очевидно, что максимум в выражении (4.99) равен 1 и достигается в точках глобальных минимумов $y = y^* \in Y^*$. При $y \notin Y^*$, «средняя» относительная плотность $\tilde{\alpha}(y, y^*) < 1$, причем её значения тем меньше, чем больше значение функции $f(y)$ отличается от глобально-оптимального. Заметим, что асимптотическое значение относительной плотности испытаний вне Y^* уменьшается с уменьшением $\delta = \lim_{k \rightarrow \infty} \delta_k$. Однако, сразу выбирать очень маленькие δ_k нельзя из-за замедления выхода относительной плотности на асимптотику. Качественный характер зависимости $\tilde{\alpha}(y, y^*)$ от y для конкретной функции показан на рис.4.6

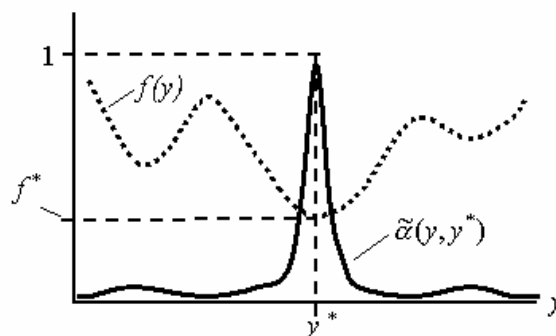


Рис.4.6. Зависимость «средней» концентрации испытаний в точке y по отношению к y^*

👉 Контрольные вопросы и упражнения.

Попробуйте самостоятельно построить оценку $\tilde{\alpha}(y, y^*)$ для метода условной глобальной оптимизации с адаптивной стохастической моделью. Рассмотрите два случая, когда точка y допустима и когда она допустимой не является. В обоих случаях определите порядки относительной концентрации испытаний.

4.6.2. Метод аналитического оценивания относительной концентрации испытаний

Перейдем теперь к формальному изложению метода получения границ колебаний относительной концентрации испытаний для некоторого подкласса T -представимых процедур.

Введем дополнительные обозначения. При $N > 1$ рассмотрим замкнутый шар

$$O(y, r) = \{u : \|y - v\| \leq r\},$$

а при $N = 1$ примем другую трактовку того же обозначения

$$O(y, r) = [y; y + r].$$

Пусть U_k — множество попарно различных точек из (4.76), и $u \in U_k$. Введем обозначение для координат точки с максимальным значением функции выбора $W_k(\omega(Q, U_k), y)$, при условии, что этот максимум определяется в области в виде шара с центром в точке испытания u и радиусом, равным значению $\bar{\rho}(u, U_k)$, т.е. расстоянию от u до ближайшей к ней другой точки из U_k .

$$\hat{y}(u, U_k) \in \text{Arg max} \{ W_k(\omega(Q, U_k), y) : y \in O(u, \bar{\rho}(u, U_k)) \}.$$

Здесь учтена возможная неединственность максимума.

Зафиксируем $u \in D, r > 0, \sigma \geq r, t \geq 2$. Рассмотрим совокупность $\psi(u, r)$, состоящую из таких множеств U_k , что $u \in U_k, \bar{\rho}(u, U_k) = r$, число точек в пересечении $O(u, \sigma) \cap U_k$ не менее t . Определим величину

$$\tilde{\rho}(u, r) = \inf \{ \rho(u, \hat{y}(u, U_k)) : U_k \in \psi(u, r) \}, \quad (4.100)$$

показывающую наименьшее возможное расстояние от существующей точки испытания u до новой, которая может быть размещена в ее окрестности, при условии, что расстояние от u до ближайшей к ней существующей точки испытания равно r и в некоторой малой окрестности u имеется достаточное количество точек испытаний.

Кроме соотношения (4.100) введем для $U_k \in \psi(u, r)$ величину

$$A(Q, U_k, u, r) = \max \{ W_k(\omega(Q, U_k), y) : y \in O(u, r) \}, \quad (4.101)$$

характеризующую максимальное значение функции выбора, соответствующее точке $\hat{y}(u, U_k)$ (заметьте, что в (4.101), из-за специфики структуры множества $\psi(u, r)$, значение $r = \bar{\rho}(u, U_k)$).

Сформулируем несколько гипотез о свойствах используемого метода поиска. Можно показать, что они выполняются для многих известных методов со всюду плотной сходимостью. Проверку справедливости этого утверждения предоставим читателю.

Пусть y' и y'' — две фиксированные точки из D .

Предположение 4.1. Для функции Q , описывающей решаемую задачу, и точек $u = y'$ и $u = y''$ из U_k существует конечная оценка ν , что

$$\nu \geq r / \tilde{\rho}(u, r) > 1 \quad (4.102)$$

при достаточно малом σ в (4.100) и $r \leq \sigma$.

Предположение 4.2. Для любой последовательности множеств U_k из (4.76) попарно различных точек, образующих в пределе всюду плотное покрытие D и включающих точки y' и y'' , найдется такое K , начиная с которого при $u = y'$,

а также для $u=y''$ при достаточно малом $\sigma \geq \bar{\rho}(u, U_k)$ существуют оценки для значений из (4.101) следующего вида

$$A_1(Q, u, \bar{\rho}(u, U_k)) \leq A(Q, U_k, u, \bar{\rho}(u, U_k)) \leq A_2(Q, u, \bar{\rho}(u, U_k)), \quad (4.103)$$

где функции $A_i(Q, u, r)$ ($i=1,2$) не зависят от U_k , σ , некоторым образом зависят от свойств Q в точке u и строго возрастают по r .

Предположение 4.3. При достаточно малом r и некотором $p>0$ существуют положительные ограниченные, непрерывные по r в нуле решения $\beta_i = \beta_i(y', y'', r)$ ($i=1,2$) системы неравенств

$$A_2(Q, y'', r) \leq A_1(Q, y', (r/\beta_1)^p) \quad (4.104)$$

$$A_2(Q, y', (r/\beta_2)^p) \leq A_1(Q, y'', r) \quad (4.105)$$

причем эти решения имеют положительные ограниченные пределы $\beta_1(y', y'', 0)$, $\beta_2(y', y'', 0)$ при $r \rightarrow 0$.

Теорема 4.13. Пусть метод поиска (4.75) на классе функций Φ обладает всюду плотным характером сходимости в области D , y' и y'' — фиксированные точки испытаний, не лежащие на границах D , а функционал $W_k(\omega(Q, Y_k), y)$ удовлетворяет предположениям 4.1–4.3. Тогда $\alpha(y', y'')$ — относительная концентрация испытаний в точке y' по отношению к y'' обладает следующими свойствами.

1. Имеет порядок p .
2. Подчиняется следующей динамике. При достаточно большем k и $r = \bar{\rho}_k(y'')$, если значение

$$\alpha_k(y', y'') \notin [\beta_1(y', y'', r); \beta_2(y', y'', r)],$$

то оно может только приближаться к этому отрезку и за конечное число шагов попадет в него. Находясь в этом отрезке $\alpha_k(y', y'')$ может как возрастать (не более, чем в $v^{1/p}$ раз на шаге), так и убывать (не более, чем в v раз на шаге).

3. Для любого $\varepsilon > 0$ найдется K , что при $k > K$ имеют место следующие оценки

$$\alpha_k(y', y'') \in [v^{-1} \beta_1(y', y'', 0) - \varepsilon; v^{1/p} \beta_2(y', y'', 0) + \varepsilon]. \quad (4.106)$$

ДОКАЗАТЕЛЬСТВО приведено в [15].

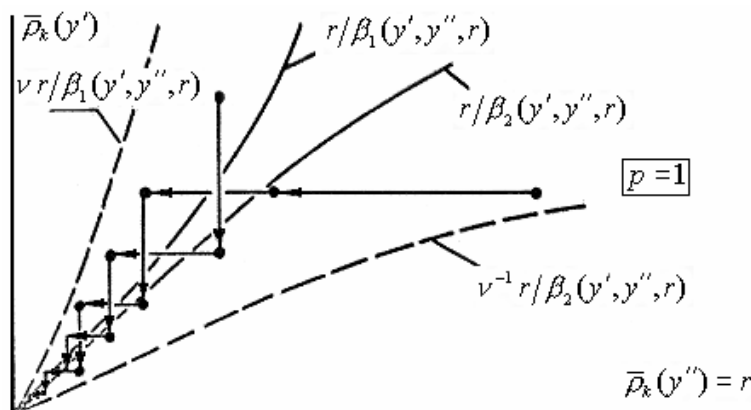


Рис.4.7. Асимптотическая динамика при $v=2$ и $p=1$

Асимптотическую динамику перемещения точки $(\bar{\rho}_k(y''); \bar{\rho}_k(y'))$ из нескольких начальных положений при $\nu=2$ и значениях $p=1$ и $p=0.5$ иллюстрируют рис.4.7, 4.8

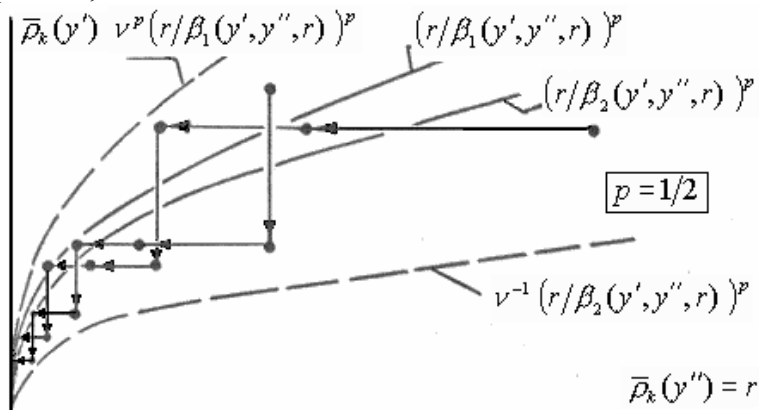


Рис.4.8. Асимптотическая динамика при $\nu=2$ и $p=0.5$

Пример оценки относительной концентрации для метода Х. Кушнера

Не проводя выкладок, укажем только окончательный результат для $y' = y$, $y'' = y^*$. Оказывается, что для метода Х. Кушнера в (4.106)

$$\beta_1(y, y^*, 0) = \beta_2(y, y^*, 0) = \tilde{\alpha}(y, y^*) = \left(\frac{1}{1 + (f(y) - f^*)/\delta} \right)^2, \quad (107)$$

и при этом $p = 1$, $\nu = 2$, т.е. при сколь угодно малом $\varepsilon > 0$, начиная с некоторого шага, относительная плотность испытаний принадлежит интервалу

$$\alpha_k(y, y^*) \in [0.5 \tilde{\alpha}(y, y^*) - \varepsilon; 2\alpha(y, y^*) + \varepsilon],$$

где $\tilde{\alpha}(y, y^*)$ определяется соотношением (107).

Пример оценки относительной концентрации испытаний для метода условной глобальной оптимизации с адаптивной вероятностной моделью

Напомним, что этот метод был описан в пункте 4.3.4, а вид его характеристической функции приведен в (4.86) раздела 4.5, где в качестве функции распределения $F(z)$ использована функция

$$F(z) = \begin{cases} 0.5 \exp(z), & z \leq 0 \\ 1 - 0.5 \exp(z), & z > 0. \end{cases}$$

Будем предполагать, что вектор-функция $Q = (f, g_1, \dots, g_m)$ принадлежит классу Φ'' , введенному при рассмотрении данного метода в разделе 4.5. Для упрощения вывода нужных нам оценок аппроксимируем все функции кусочно-постоянными из Φ'' , т.е. рассмотрим подкласс в Φ'' .

Опуская достаточно громоздкие, но несложные выкладки, приведем окончательный результат. Рассмотрим подобласть Y_{ε^G} допустимой области Y из (4.38)

$$Y_{\varepsilon^G} = \{y \in D : G(y) \leq -\varepsilon^G\},$$

где $0 < \varepsilon^G \ll 1$, а $G(y)$ — обобщенная функция ограничения, введенная в п. 4.3.5. Можно показать, что если обе сравниваемые точки допустимы с «запасом»,

т.е. $y \in Y_{\varepsilon^g}$ и $y^* \in Y_{\varepsilon^g}$, и точка $y^* \in Y^*$, то порядок p относительной плотности $\alpha_k(y, y^*)$ испытаний в допустимой точке y по отношению к точке глобального минимума равен 1, причем в выражении (4.106) для границ колебаний относительной концентрации $\nu = 2$, а значения

$$\beta_1(y, y^*, 0) = 0.5 \tilde{\alpha}(y, y^*), \quad \beta_2(y, y^*, 0) = 2 \tilde{\alpha}(y, y^*),$$

$$\tilde{\alpha}(y, y^*) = \left(\frac{1}{1 + (f(y) - f^*)/\delta} \right)^{1/s}, \quad (108)$$

т.е. при сколь угодно малом $\varepsilon > 0$, начиная с некоторого шага, относительная концентрация испытаний принадлежит интервалу

$$\alpha_k(y, y^*) \in [0.25 \tilde{\alpha}(y, y^*) - \varepsilon; 4 \alpha(y, y^*) + \varepsilon],$$

где $\tilde{\alpha}(y, y^*)$ определяется соотношением (108).

Если же $y^* \in Y_{\varepsilon^g} \cap Y^*$, а точка $y \notin Y$, то порядок относительной концентрации испытаний в недопустимой точке y по отношению к точке глобального минимума определяется как $p = (s/\gamma) < 1$. Т.е. при $k \rightarrow \infty$ расстояние $\bar{\rho}_k(y^*)$ от y^* до ближайшей к ней точки испытания является величиной более высокого порядка малости (равного $1/p > 1$) нежели расстояние $\bar{\rho}_k(y)$ от точки y недопустимого испытания до ближайшей к ней точки другого испытания:

$$\bar{\rho}_k(y^*) = o(\bar{\rho}_k(y)) = O((\bar{\rho}_k(y))^{1/p}).$$

Кроме того, в выражении (4.106) для границ колебаний относительной концентрации $\nu = 2$, а значения

$$\beta_1(y, y^*, 0) = 2^{-1/p} \tilde{\alpha}(y, y^*), \quad \beta_2(y, y^*, 0) = 2 \tilde{\alpha}(y, y^*),$$

$$\tilde{\alpha}(y, y^*) = \left(\frac{\delta}{G(y) + \varepsilon_g} \right)^{1/s}, \quad (109)$$

т.е. при сколь угодно малом $\varepsilon > 0$, начиная с некоторого шага, относительная концентрация испытаний принадлежит интервалу

$$\alpha_k(y, y^*) \in [2^{-(1+\gamma/s)} \tilde{\alpha}(y, y^*) - \varepsilon; 2^{(1+\gamma/s)} \alpha(y, y^*) + \varepsilon],$$

где $\tilde{\alpha}(y, y^*)$ определяется соотношением (109).

Таким образом, для рассмотренного метода условной оптимизации, обладающего всюду плотным характером сходимости, концентрация точек испытаний в недопустимой области имеет более высокий порядок малости, по сравнению с их концентрацией в допустимой области в окрестности решения. Кроме того, обеспечивается значительно большая плотность испытаний в окрестностях решений, нежели в других допустимых точках. Полученные соотношения дают количественное описание динамики относительной концентрации испытаний.

Материал раздела 4.6 позволяет увидеть, что методы многоэкстремальной оптимизации со всюду плотным в пределе характером размещения испытаний в действительности могут являться эффективным средством решения задач многоэкстремальной оптимизации за счет способности к построению существенно неравномерных покрытий, адаптирующихся к решаемой задаче.

Лист регистрации изменений

Дата	Автор	Комментарии
??.06.03	Гришагин В.А.	Первоначальная версия главы 4
??.07.03	Гришагин В.А.	Правка текста
13.05.03	Городецкий С.Ю.	Создание копии раздела 4.5 с использованием версии раздела 4.5, написанной В.А.Гришагиным
14.08.03	Городецкий С.Ю.	Начата переработка раздела 4.5
24.08.03	Городецкий С.Ю.	Закончена переработка раздела 4.6
24.08.03	Городецкий С.Ю.	Начато добавление раздела 4.6
25.08.03	Городецкий С.Ю.	Закончен раздел 4.6 + рисунки
25.09.03	Городецкий С.Ю.	Внесение изменений в раздел 4.5
25.09.03	Городецкий С.Ю.	Корректурa разделов 4.5, 4.6
25.09.03	Городецкий С.Ю.	Внесение изменений в раздел 4.6
09.10.03	Городецкий С.Ю.	Дополнения в раздел 4.5
11.10.03	Городецкий С.Ю.	Окончательная корректурa 4.5–4.6
09.10.03– 17.10.03	Гришагин В.А.	Изменение обозначений в разделах 4.1–4.4, корректурa
18.10.03	Городецкий С.Ю.	Изменение стилей оформления в разделах 4.1–4.4
18.10.03	Городецкий С.Ю.	Присоединение новых версий разделов 4.5–4.6 к главе 4

Глава 5. Фундаментальные способы редукции размерности в многоэкстремальных задачах

5.1. Многоэкстремальные задачи и методы покрытий

Рассмотрим конечномерную задачу оптимизации (задачу нелинейного программирования)

$$f(y) \rightarrow \inf, y \in Y \subseteq R^N \quad (5.1)$$

$$Y = \{y \in D : g_j(y) \leq g_j^+, 1 \leq j \leq m\} \quad (5.2)$$

$$D = \{y \in R^N : y_i \in [a_i, b_i], 1 \leq i \leq N\}, \quad (5.3)$$

т.е. задачу отыскания экстремальных значений целевой (минимизируемой) функции $f(y)$ в области Y , задаваемой координатными (5.3) и функциональными (5.2) ограничениями на выбор допустимых точек (векторов) $y = (y_1, y_2, \dots, y_N)$. В данной модели допуски g_j^+ , ограничивающие сверху допустимые изменения функций $g_j(y)$, $1 \leq j \leq m$, являются константами, а величины a_i, b_i , $1 \leq i \leq N$, задающие границы изменения варьируемых параметров задачи (координат вектора y) либо константы, либо, когда соответствующая нижняя и (или) верхняя границы отсутствуют, принимаются равными $a_i = -\infty$ и (или) $b_i = +\infty$.

Довольно часто в формулировку задачи нелинейного программирования включают также ограничения в виде равенств. Однако любое равенство $h(y) = 0$, во-первых, формально можно представить в виде системы двух неравенств $h(y) \leq 0$ и $-h(y) \leq 0$. Во-вторых, при численном решении задачи оптимизации на ЭВМ точная реализуемость равенства невозможна, поэтому предполагают, что допустима его выполнимость с некоторой погрешностью $\delta > 0$, т.е. вместо равенства $h(y) = 0$ рассматривается неравенство $|h(y)| \leq \delta$. Таким образом, учитывая указанные обстоятельства, можно утверждать, что (5.1)-(5.3) является формулировкой общей многомерной задачи нелинейного программирования.

Если $m = 0$, т.е. функциональные ограничения отсутствуют, будем полагать $Y = D$. Задача (5.1)-(5.3) в этом случае будет называться задачей безусловной оптимизации.

Предметом рассмотрения настоящей главы являются многоэкстремальные задачи оптимизации, т.е. задачи, в которых целевая функция $f(y)$ имеет в допустимой области Y несколько локальных экстремумов. На сложность решения таких задач существенное влияние оказывает размерность. Например, для класса многоэкстремальных функций, удовлетворяющих условию Липшица, имеет место так называемое "проклятие размерности", состоящее в экспоненциальном росте вычислительных затрат при увеличении размерности. А именно: если в одномерной задаче для достижения точности решения ε требуется p вычислений функции, то в задаче с размерностью N для решения с той же точностью необходимо осуществить αp^N испытаний, где α зависит от целевой функции, допустимой области и используемого метода.

Для специальных узких классов многоэкстремальных задач порядок роста затрат может быть и лучше экспоненциального. Например, как мы увидим далее, для сепарабельных функций, минимизируемых в гиперинтервале D , затраты растут линейно. Этот факт в очередной раз иллюстрирует то обстоятельство, что улучшить эффективность решения можно только на основе глубокого учета априорной информации о задаче.

Формой такого учета может быть построение эффективных методов оптимизации как оптимальных решающих правил на основе минимаксного или байесовского подходов к понятию оптимальности метода поиска экстремума в рамках соответствующих математических моделей. К сожалению, для многомерных многоэкстремальных задач проблема построения оптимальных алгоритмов является очень сложной, и построить оптимальные (в том или ином смысле алгоритмы) удастся в исключительных случаях.

В силу того, что глобальный экстремум является не локальной, а интегральной характеристикой минимизируемой функции, т.е. для его определения необходимо сопоставить значение в точке глобального минимума со значениями функции в остальных точках области поиска, поисковый метод многоэкстремальной оптимизации вынужден строить некоторую сетку испытаний, покрывающую допустимую область. Свойства этой сетки будут определять эффективность метода и в свою очередь формироваться как следствие свойств модели задачи, на основе которой построен применяемый алгоритм.

Так, если модель задачи ориентирована на класс непрерывных функций, то, как показано в предыдущей главе, построение оценок оптимума по конечному числу испытаний невозможно, и единственный способ достижения решения состоит в построении всюду плотной последовательности (синонимично – сетки) испытаний.

Другой пример. Пусть минимизируемая функция $f(y)$ удовлетворяет в области Y условию Липшица (4.69) с константой L , или кратко $f(y) \in Lip(Y)$. В этом случае оказывается, что задача построения эффективной конечной сетки эквивалентна геометрической задаче построения оптимального покрытия области поиска [11, 45]. При ориентации на критерий априорной оптимальности (4.26) такими покрытиями являются так называемые равномерные, или переборные сетки. Из них наиболее элементарными являются кубические сетки. Методика их построения чрезвычайно проста.

Рассмотрим единичный N -мерный гиперкуб $D = \{y \in \mathbb{R}^N : y_i \in [0, 1], 1 \leq i \leq N\}$. Выберем целое число $M > 0$, разделим каждое ребро куба, лежащее на одной из координатных осей, на M равных частей и в точках деления построим плоскости, перпендикулярные ребру. Полученные плоскости разделят гиперкуб на M^N равных кубических элементов, центры которых составят кубическую сетку. На первый взгляд, точки этой сетки расположены в гиперкубе с хорошей "равномерностью". Однако, подобное утверждение справедливо лишь в одномерном случае. Уже при $N = 2$ кубическая решетка не очень хороша, а с увеличением размерности ее "равномерность" быстро ухудшается. Существуют сетки, реализующие более качественное покрытие области D в том смысле, что они обеспечивают меньшие значения показателей "неравномерности" размещения узлов сетки (например, сетки метода Монте-Карло). Среди этих сеток наименьшим порядком роста неравномерности при увеличении размерности обладают так называемые ЛП _{τ} -сетки [41] (теоретически наилучший порядок роста неравномерности неизвестен).

Размещение испытаний в узлах кубических, или ЛП_r-сеток, или любых других с заранее известными узлами и последующий выбор наименьшего вычисленного значения функции как оценки экстремума порождает класс пассивных, или *неадаптивных* алгоритмов оптимизации. При наличии аппарата оценивания погрешности решения существенного увеличения эффективности можно достичь, если выбирать узлы сетки на основе полученной информации об исследуемой функции. Такие сетки называются *адаптивными*, и некоторые способы их построения мы будем рассматривать в дальнейшем.

5.2. Принципы редукции сложности в многомерных многоэкстремальных задачах

Другой подход к конструированию численных методов анализа многомерных многоэкстремальных задач использует идею *редукции сложности*, когда решение исходной задачи заменяется решением одной или нескольких более простых задач.

Одна из таких схем редукции базируется на том элементарном факте, что глобальный экстремум является локальным. Отсюда следует прозрачный вывод, что для нахождения глобально-оптимального решения достаточно найти все локальные минимумы и из них выбрать наименьший. В контексте теоретического подхода эта конструкция безупречна, однако, на практике не все так безоблачно.

Прежде всего, возникает вопрос: как найти все локальные минимумы? Эта задача может быть решена, если для каждого локального минимума известна зона его притяжения, т.е. такая окрестность точки минимума, в которой функция унимодальна. Тогда, поместив начальную точку поиска в эту окрестность, можно тем или иным локальным методом найти искомым локальный минимум. Другими словами, данная схема предполагает известным разбиение области поиска на зоны притяжения локальных минимумов. Однако на практике подобная априорная информация, как правило, отсутствует (даже количество локальных экстремумов обычно не известно). Поэтому при таком подходе возникает дополнительная задача выбора начальных точек.

Простейший способ состоит в том, чтобы выбирать начальные точки по схеме метода Монте-Карло [27], т.е. случайно в соответствии с некоторым распределением в области поиска (обычно равномерным), либо использовать в качестве таких точек узлы некоторой регулярной сетки [41]. Сравнение способов выбора начальных точек приведено в [45].

При таком способе сходимость к глобальному экстремуму обеспечивается тем обстоятельством, что при увеличении числа начальных точек хотя бы одна из них попадет в зону притяжения глобального экстремума, так как эти сетки (регулярная, либо случайная) строятся так, чтобы обеспечить уплотняющееся покрытие всей области поиска.

У этой схемы есть, однако, существенный недостаток. Дело в том, что в зону притяжения одного и того же локального минимума могут попасть несколько начальных точек, т.е. придется несколько раз искать один и тот же локальный минимум. Для устранения этого недостатка предложено несколько различных схем.

Идея одной из них состоит в следующем. Вначале в области поиска выбираются L базовых точек (обычно более или менее равномерно расположенных в области поиска) и затем вычисляются значения функции в этих точках, т.е., по существу, производится грубая оценка поведения функции. Из

множества базовых точек производится отбор $l < L$ "лучших" точек, т.е. таких, в которых значения функции минимальны и точки не слишком близки друг к другу. Эти l точек и служат начальными точками локального поиска. Существует множество вариантов этой схемы. Например, семейство базовых точек может модифицироваться в процессе поиска – туда могут добавляться новые точки. Также в процессе реализации этой схемы может изменяться схема отбора точек с учетом вновь поступившей информации и т.п.

Существуют и другие схемы, построенные на основе редукции многоэкстремальной задачи к решению локальных подзадач. Алгоритмы, сконструированные в рамках данного подхода, обладают асимптотической сходимостью к глобальному экстремуму при слабых предположениях о непрерывности или той или иной степени дифференцируемости целевой функции. Эта сходимость обеспечивается тем свойством алгоритмов, что любая точка области поиска является предельной точкой последовательности испытаний $\{y^k\}$, порождаемой алгоритмом, и, следовательно, для непрерывной функции

$$\lim_{\tau \rightarrow \infty} \min_{1 \leq k \leq \tau} f(y^k) = \min_{y \in Y} f(y) \quad (5.4)$$

Другими словами, последовательность испытаний этих методов является всюду плотной в области поиска, т.е. сходится в смысле Определения 4.2 ко всем точкам этой области, а, следовательно, и к точке глобального минимума.

Как отмечено в предыдущей главе, указанный характер сходимости является слишком "расточительным", что не способствует эффективности методов данного типа. Избежать указанного недостатка можно только при наличии достаточной и довольно богатой априорной информации о задаче (типа сведений об областях притяжения). Такая информация довольно редко имеет место на практике и, как правило, может быть получена только для простых многоэкстремальных задач с небольшим числом экстремумов, а для существенно многоэкстремальных задач применение идеи сведения к задачам на локальный минимум не слишком эффективно.

По сравнению с редукцией к локальным задачам более плодотворным является подход, основанный на идеях *редукции размерности*, т.е. построение таких схем, когда решение многомерной задачи сводится к решению одной или нескольких одномерных подзадач. Известны две эффективные схемы такого рода: многошаговая схема оптимизации [7, 9, 42] и редукция размерности на основе кривых, заполняющих пространство (кривых Пеано) [7, 9].

Многошаговая схема редукции размерности сводит решение задачи (5.1)-(5.3) к решению семейства рекурсивно связанных одномерных подзадач, в которых каждое вычисление целевой функции одномерной подзадачи есть решение новой одномерной подзадачи следующего уровня рекурсии, за исключением последнего уровня, где вычисляется значение исходной многомерной целевой функции.

Вторая из указанных схем основана на известном фундаментальном факте, согласно которому N -мерный гиперпараллелепипед (5.3) и отрезок $[0, 1]$ вещественной оси являются равномошными множествами и отрезок $[0, 1]$ может быть однозначно и непрерывно отображен на гиперпараллелепипед (5.3) [29]. Отображения такого рода обычно называют *кривыми*, или *развертками Пеано*.

Пусть $y(x), x \in [0, 1]$, есть кривая Пеано и функция $f(y)$ из (5.1) непрерывна. Тогда из непрерывности $f(y)$ и $y(x)$ и равенства

$$D = \{y(x) : 0 \leq x \leq 1\}$$

следует, что

$$\min_{y \in D} f(y) = \min_{x \in [0,1]} f(y(x)),$$

т.е. решение многомерной задачи минимизации $f(y)$ сводится к минимизации одномерной функции $f(y(x))$.

Рассмотрение многошаговой схемы оптимизации и схем редукции размерности на основе кривых Пеано станет основным содержанием настоящей главы.

5.3. Многошаговая схема редукции размерности

Рассмотрим в качестве исходной задачу (5.1)-(5.3) в варианте, когда все допуски $g_j^+ = 0$, $1 \leq j \leq m$. Это несколько не ограничивает общности анализа (всегда вместо ограничений $g_j(y) \leq g_j^+$ можно ввести ограничения $\tilde{g}(y) = g_j(y) - g_j^+ \leq 0$), однако, позволяет упростить изложение.

Предположим также, что все функции-ограничения $g_j(y)$, $1 \leq j \leq m$, являются непрерывными, а область D - ограниченной, что обеспечивает компактность области Y .

Введем непрерывную функцию, определенную в области D , такую, что

$$\begin{aligned} G(y) &\leq 0, y \in Y \\ G(y) &> 0, y \notin Y \end{aligned} \quad (5.5)$$

В качестве $G(y)$ можно взять, например,

$$G(y) = \max\{g_1(y), \dots, g_m(y)\} \quad (5.6)$$

или

$$G(y) = \max\{0; g_1(y), \dots, g_m(y)\} \quad (5.7)$$

Последняя функция тождественно равна нулю в области Y .

Введем обозначения

$$u_i = (y_1, \dots, y_i), v_i = (y_{i+1}, \dots, y_N), \quad (5.8)$$

позволяющие при $1 \leq i \leq N-1$ записать вектор y в виде пары $y = (u_i, v_i)$, и примем, что $y = v_0$ при $i = 0$ и $y = u_N$ при $i = N$.

Введем сечения множества Y :

$$S_1 = Y, \quad S_{i+1}(u_i) = \{(u_i, v_i) \in Y\}, \quad 1 \leq i \leq N-1, \quad (5.9)$$

и проекции сечений на ось y_{i+1} :

$$\Pi_{i+1}(u_i) = \{y_{i+1} \in R^1 : \exists (y_{i+1}, v_{i+1}) \in S_{i+1}(u_i)\}, \quad (5.10)$$

Положим $G^N(y) \equiv G(y)$ и построим семейство функций

$$G^i(u_i) = \min\{G^{i+1}(u_i, y_{i+1}) : y_{i+1} \in [a_{i+1}, b_{i+1}]\}, \quad 1 \leq i \leq N-1, \quad (5.11)$$

определенных в соответствующих проекциях

$$D_i = \{u_i \in R^i : y_j \in [a_j, b_j], 1 \leq j \leq i\} \quad (5.12)$$

множества D из (5.3) на координатные оси y_1, \dots, y_i , причем по определению $D_N = D$.

В силу непрерывности функции $G^N(y) \equiv G(y)$ и компактности области D функция $G^{N-1}(u_{N-1})$ существует и непрерывна в D_{N-1} , что влечет существование и непрерывность функции $G^{N-2}(u_{N-2})$ и т.д. определяет существование и непрерывность всех функций семейства (5.11).

Справедлива следующая

Лемма 5.1. *Справедливо соотношение*

$$G^i(u_i) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\}, 1 \leq i \leq N-1. \quad (5.13)$$

ДОКАЗАТЕЛЬСТВО. Вследствие (5.11) доказательство (5.13) состоит в установлении справедливости равенства

$$\min_{y_{i+1} \in [a_{i+1}, b_{i+1}]} \dots \min_{y_N \in [a_N, b_N]} G(u_i, v_i) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\} \quad (5.14)$$

В силу непрерывности $G(y)$ для любого $u_i \in D_i$ существует $\bar{v}_i = (\bar{y}_{i+1}, \dots, \bar{y}_N)$ такой, что $\bar{y}_j \in [a_j, b_j]$, $i+1 \leq j \leq N$, и

$$\min_{y_{i+1} \in [a_{i+1}, b_{i+1}]} \dots \min_{y_N \in [a_N, b_N]} G(u_i, v_i) = G(u_i, \bar{v}_i),$$

откуда следует, что левая часть равенства (5.14) больше или равна правой.

Покажем обратное неравенство. В силу непрерывности $G(y)$ и компактности D существует вектор $v_i^* = (y_{i+1}^*, \dots, y_N^*)$ такой, что

$$G(u_i, v_i^*) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\}.$$

Согласно (5.11) имеем:

$$G^{N-1}(u_i, y_{i+1}^*, \dots, y_{N-1}^*) = \min\{G(u_i, y_{i+1}^*, \dots, y_{N-1}^*, y_N) : y_N \in [a_N, b_N]\} \leq G(u_i, v_i^*),$$

$$G^{N-2}(u_i, y_{i+1}^*, \dots, y_{N-2}^*) = \min\{G^{N-1}(u_i, y_{i+1}^*, \dots, y_{N-2}^*, y_{N-1}) : y_{N-1} \in [a_{N-1}, b_{N-1}]\} \leq \\ \leq G^{N-1}(u_i, y_{i+1}^*, \dots, y_{N-1}^*),$$

...

$$G^i(u_i) = \min\{G^{i+1}(u_i, y_{i+1}) : y_{i+1} \in [a_{i+1}, b_{i+1}]\} \leq G^{i+1}(u_i, y_{i+1}^*) \leq \dots \leq G(u_i, v_i^*).$$

Лемма доказана.

Введем проекции

$$Y_i = \{u_i \in R^i : \exists (u_i, v_i) \in Y\}, 1 \leq i \leq N, \quad (5.15)$$

множества Y на координатные оси y_1, \dots, y_i .

Лемма 5.2. *Представление (5.15) эквивалентно соотношению*

$$Y_i = \{u_i \in R^i : G^i(u_i) \leq 0\} \quad (5.16)$$

ДОКАЗАТЕЛЬСТВО. Пусть выполнено (5.15), т.е. для некоторого u_i существует v_i^* такой, что $(u_i, v_i^*) \in Y$. Но тогда $G(u_i, v_i^*) \leq 0$, т.е. вследствие Леммы 5.1

$$G^i(u_i) = \min \{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\} \leq G(u_i, v_i^*) \leq 0,$$

и (5.16) справедливо.

Пусть теперь, наоборот, для некоторого $u_i \in Y_i$ выполняется $G^i(u_i) \leq 0$. Но согласно Лемме 5.1 существует вектор v_i^* такой, что $G^i(u_i) = G(u_i, v_i^*)$, т.е. $(u_i, v_i^*) \in Y$.

Лемма доказана.

Лемма 5.3. *Определение проекции (5.10) эквивалентно представлению*

$$\Pi_{i+1}(u_i) = \{y_{i+1} \in [a_{i+1}, b_{i+1}] : G^{i+1}(u_i, y_{i+1}) \leq 0\}, \quad (5.17)$$

ДОКАЗАТЕЛЬСТВО. Прежде всего заметим, что поскольку $Y \subseteq D$, то необходимо $a_j \leq y_j \leq b_j, 1 \leq j \leq N$.

Пусть теперь y_{i+1} таков, что $G^{i+1}(u_i, y_{i+1}) \leq 0$. Тогда существует v_{i+1}^* такой, что $G(u_i, y_{i+1}, v_{i+1}^*) \leq 0$, т.е. $(y_{i+1}, v_{i+1}^*) \in S_{i+1}(u_i)$, следовательно, y_{i+1} принадлежит проекции $\Pi_{i+1}(u_i)$ в смысле определения (5.10).

Предположим далее, что для некоторого y_{i+1} существует v_{i+1}^* такой, что $(y_{i+1}, v_{i+1}^*) \in S_{i+1}(u_i)$. Тогда вектор $(u_i, y_{i+1}, v_{i+1}^*) \in Y$, т.е. $G^{i+1}(u_i, y_{i+1}) \leq G(u_i, y_{i+1}, v_{i+1}^*) \leq 0$.

Лемма доказана.

 **Контрольные вопросы и упражнения:**

1. Докажите, что сечение $S_{i+1}(u_i)$ не пусто тогда и только тогда, когда $G^i(u_i) \leq 0$.
2. Докажите, что неравенство $G^i(u_i) \leq 0$ является необходимым и достаточным условием непустоты проекции $\Pi_{i+1}(u_i)$.

Предположим теперь непрерывность функции $f(y)$ и, положив по определению $f^N(y) \equiv f(y)$, построим семейство функций

$$f^i(u_i) = \min \{f^{i+1}(u_i, y_{i+1}) : y_{i+1} \in \Pi_{i+1}(u_i)\}, \quad 1 \leq i \leq N-1, \quad (5.18)$$

определенных на соответствующих проекциях Y_i . Тогда имеет место основное соотношение

$$\min_{y \in Y} f(y) = \min_{y_1 \in \Pi_1} \min_{y_2 \in \Pi_2(u_1)} \dots \min_{y_N \in \Pi_N(u_{N-1})} f(y) \quad (5.19)$$

Как следует из (5.19), для решения задачи (5.1) – (5.3) достаточно решить одномерную задачу

$$f^1(y_1) \rightarrow \min_{y_1 \in \Pi_1 \subseteq R^1} \quad (5.20)$$

$$\Pi_1 = \{y_1 \in [a_1, b_1] : G^1(y_1) \leq 0\} \quad (5.21)$$

При этом каждое вычисление функции $f^1(y_1)$ в некоторой фиксированной точке $y_1 \in \Pi_1$ представляет собой согласно (5.18) решение одномерной задачи

$$f^2(y_1, y_2) \rightarrow \min_{y_2 \in \Pi_2(y_1) \subseteq R^1} \quad (5.22)$$

$$\Pi_2(y_1) = \{y_2 \in [a_2, b_2] : G^2(y_1, y_2) \leq 0\} \quad (5.23)$$

Эта задача является одномерной задачей минимизации по y_2 , т.к. y_1 фиксировано.

В свою очередь, каждое вычисление значения функции $f^2(y_1, y_2)$ при фиксированных y_1, y_2 требует решения одномерной задачи

$$f^3(u_2, y_3) \rightarrow \min_{y_3 \in \Pi_3(u_2)} \quad (5.24)$$

и т.д. вплоть до решения задачи

$$f^N(u_{N-1}, y_N) = f(u_{N-1}, y_N) \rightarrow \min_{y_N \in \Pi_N(u_{N-1})} \quad (5.25)$$

при фиксированном u_{N-1} .

Окончательно решение задачи (5.1) – (5.3) сводится к решению семейства "вложенных" одномерных подзадач

$$f^i(u_{i-1}, y_i) \rightarrow \min_{y_i \in \Pi_i(u_{i-1})}, \quad (5.26)$$

где фиксированный вектор $u_{i-1} \in Y_{i-1}$.

Решение исходной многомерной задачи (5.1) – (5.3) через решение системы взаимосвязанных одномерных подзадач (5.26) называется *многошаговой схемой редукции размерности*.

Заметим, что если в задачах (5.26) по некоторым координатам y_i искать не минимум, а максимум, то становится возможным вычисление сложных минимаксных (максиминных) выражений вида

$$\underset{y_1 \in \Pi_1}{extr} \underset{y_2 \in \Pi_2(u_1)}{extr} \dots \underset{y_N \in \Pi_N(u_{N-1})}{extr} f(y)$$

где операция *extr* означает вычисление глобального минимума или максимума.

Проиллюстрируем общие результаты конкретным несложным примером.

Пример 5.1. Рассмотрим двумерную задачу (5.1) – (5.3), в которой

$$f(y) = y_1^2 + y_2^2, \quad Y = \{y \in D : (y_1 - 1)^2 + (y_2 - 1)^2 - 1 \leq 0\}, \quad (5.27)$$

$$D = \{y \in R^2 : 0 \leq y_1, y_2 \leq 2\}.$$

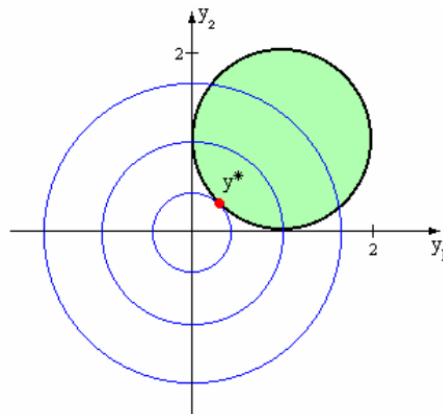


Рис. 5.1. Иллюстрация к примеру на метод редукции по многошаговой схеме

На рис. 5.1 цветом отмечена допустимая область, а концентрические круги показывают линии уровня целевой функции.

В этой задаче функция $G^2(y) = (y_1 - 1)^2 + (y_2 - 1)^2 - 1$, а функция

$$G^1(y) = \min\{G^2(y_1, y_2 : y_2 \in [0, 2])\} = (y_1 - 1)^2 - 1,$$

поскольку функция $G^2(y)$ достигает своего минимума по y_2 на отрезке $[0, 2]$ в точке $y_2 = 1$.

Согласно (5.17) области неположительности функций $G^1(y_1)$ по y_1 и $G^2(y_1, y_2)$ по y_2 определяют проекции Π_1 и $\Pi_2(y_1)$ соответственно. Границы областей задаются корнями указанных функций, принадлежащими отрезку $[0, 2]$. Для функции $G^1(y_1)$ такими корнями являются значения 0 и 2, поэтому $\Pi_1 = [0, 2]$.

Функция $G^2(y) = (y_2 - 1)^2 - \alpha^2$, где $\alpha = \sqrt{1 - (y_1 - 1)^2} \leq 1$, имеет корни $1 \pm \alpha$, очевидно принадлежащие отрезку $[0, 2]$, и неположительна между этими корнями, поэтому

$$\Pi_2(y_1) = [1 - \sqrt{1 - (y_1 - 1)^2}, 1 + \sqrt{1 - (y_1 - 1)^2}] \quad (5.28)$$

Функция $f(y) = y_1^2 + y_2^2$, будучи возрастающей по y_2 на отрезке $[0, 2]$ и, следовательно, в области (5.28), достигает своего минимума в точке $1 - \alpha$, поэтому

$$f^1(y_1) = y_1^2 + (1 - \sqrt{1 - (y_1 - 1)^2})^2 = 1 + 2y_1 - 2\sqrt{1 - (y_1 - 1)^2}.$$

Первая производная $(f(y_1))' = 2 + \frac{2(y_1 - 1)}{\sqrt{1 - (y_1 - 1)^2}}$ имеет единственный

корень $y_1^* = 1 - \frac{1}{\sqrt{2}}$. К тому же вторая производная $(f(y_1))'' = \frac{2}{(1 - (y_1 - 1)^2)^{3/2}} > 0$,

поэтому точка $y_1^* = 1 - \frac{1}{\sqrt{2}}$ доставляет минимальное значение функции $f^1(y_1)$,

равное $3 - 2\sqrt{2}$. Это значение и является искомым минимальным значением функции $f(y)$ в области Y . Для определения координаты y_2 , которая совместно

с $y_1^* = 1 - \frac{1}{\sqrt{2}}$ задает точку минимума, рассмотрим функцию

$f^2(y_1^*, y_2) = \left(1 - \frac{1}{\sqrt{2}}\right)^2 + y_2^2$ и найдем ее минимум в области

$\Pi_2(y_1^*) = \left[1 - \frac{1}{\sqrt{2}}, 1 + \frac{1}{\sqrt{2}}\right]$, который очевидно достигается в точке $y_2^* = 1 - \frac{1}{\sqrt{2}}$.

Как итог, решением исходной многомерной задачи (5.27) является вектор $y^* = \left(1 - \frac{1}{\sqrt{2}}, 1 - \frac{1}{\sqrt{2}}\right)$, обеспечивающий минимальное значение целевой функции $f(y^*) = 3 - 2\sqrt{2}$. На рисунке координата оптимума помечена точкой на соприкосновении границы допустимой области с одной из линий уровня.

В рассмотренном примере мы построили границы областей одномерного поиска аналитически, установив области неположительности соответствующих функций $G^i(u_i)$. Вместе с тем можно указать более наглядный "геометрический" способ построения проекций $\Pi_{i+1}(u_i)$. Собственно, этот способ вытекает из определений (5.9) и (5.10) и состоит в том, что необходимо построить сечения области Y плоскостями $u_i = const$ и затем установить границы этих сечений по координате y_{i+1} .

В связи с этим обратим внимание на следующее. Вычисление глобального минимума в соответствии с соотношением (5.19) аналогично процедуре нахождения многомерного интеграла от функции $f(y)$ в области Y посредством сведения к вычислению повторных одномерных интегралов. При этом области одномерного интегрирования как раз и являются соответствующими проекциями $\Pi_{i+1}(u_i)$.

Для иллюстрации рассмотрим следующий

Пример 5.2. Пусть область оптимизации (или интегрирования) задается как

$$Y = \{y \in R^2 : -4 \leq y_1, y_2 \leq 4, y_1^2 + y_2^2 \leq 4, y_1^2 + y_2^2 \geq 1\} \quad (5.29)$$

Ее вид показан на рис. 5.2.

Серый цвет, как и ранее, помечает допустимую область. Сплошная прямая $y_1 = const$ на пересечении с допустимой областью формирует сечение $S_1(y_1)$, а его проектирование на ось y_2 определяет проекцию $\Pi_2(y_1)$. Данная проекция в виде двух отрезков на оси y_2 , а также проекция Π_1 на оси y_1 отмечены на рисунке утолщенными линиями.

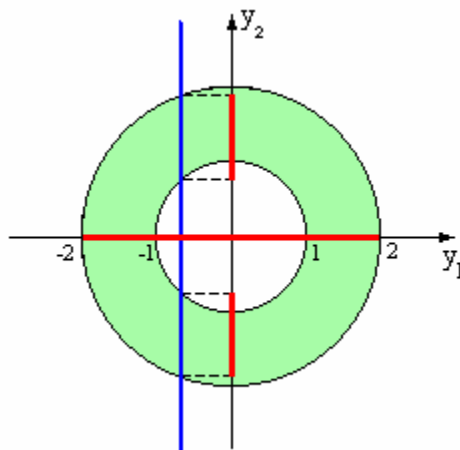


Рис. 5.2. Иллюстрация к построению проекций

Подобные геометрические соображения позволяют построить требуемые проекции задачи как

$$\Pi_1 = [-2, 2],$$

$$\Pi_2(y_1) = \begin{cases} [-\sqrt{4-y_1^2}, \sqrt{4-y_1^2}], y_1 \in [-2, -1] \cup [1, 2]; \\ [-\sqrt{4-y_1^2}, -\sqrt{1-y_1^2}] \cup [\sqrt{1-y_1^2}, \sqrt{4-y_1^2}], y_1 \in [-1, 1]. \end{cases}$$

 **Контрольные вопросы и упражнения:**

1. Постройте проекции Π_1 и $\Pi_2(y_1)$ для области $Y = \tilde{Q} \cup \hat{Q}$, где

$$\tilde{Q} = \{y \in R^2 : -0.5 \leq y_1 \leq 1.5, -1 \leq y_2 \leq 1, |y_1| + |y_2| \leq 1\},$$

$$\hat{Q} = \{y \in R^2 : (y_1 - 6)^2 + (y_2 - 4)^2 \leq 4\}.$$

2. Решите аналитически задачу минимизации функции $f(y) = y_1^2 + y_2^2$ в области (5.29).

Рассмотренные примеры являются достаточно простыми в том смысле, что нам удалось в явном виде выписать границы областей одномерного поиска – проекций Π_i – и аналитически решить одномерные задачи (5.18). Реальные практические задачи, разумеется, гораздо сложнее и не поддаются аналитическому решению. В чем же состоит эта сложность?

Обратим внимание, что при анализе одномерных подзадач многошаговой схемы возникают две проблемы:

а) необходимо сконструировать допустимые области одномерного поиска $\Pi_i(u_{i-1})$;

б) требуется обеспечить минимизацию одномерных функций $f^i(u_{i-1}, y_i)$ в областях $\Pi_i(u_{i-1})$.

Структура и сложность проекций $\Pi_i(u_{i-1})$ полностью определяются сложностью многомерной допустимой области Y . Сложность второй проблемы зависит от характеристик функций $f^i(u_i)$, на которые влияют как свойства целевой функции $f(y)$, так и особенности области поиска Y , определяемые ограничениями (5.2), (5.3).

5.4. Свойства одномерных подзадач многошаговой схемы

5.4.1. Структура допустимых областей одномерного поиска

Для анализа структуры областей $\Pi_i(u_{i-1})$ используем результаты Леммы 5.3, которая установила эквивалентность определения (5.10) и представления (5.17). Дело в том, что (5.17) дает конструктивный аппарат построения области $\Pi_i(u_{i-1})$, связывая ее с областью неположительности функции $G^i(u_{i-1}, y_i)$.

Т.к. функция $G(y)$ предполагается непрерывной в области D , то функции $G^i(u_i)$ также являются непрерывными по $u_i \in D_i$ из (5.12), а тем самым и по $y_i \in [a_i, b_i]$. Тогда при фиксированном $u_{i-1} \in D_{i-1}$ каждая из одномерных задач (5.26) является задачей вида

$$\varphi(x) \rightarrow \min, x \in \bar{Q} \subset R^1, \quad (5.30)$$

$$\bar{Q} = \{x \in [a, b] : g(x) \leq 0\},$$

причем функция $g(x)$ непрерывна.

Непрерывность ограничения $g(x)$ позволяет утверждать, что допустимая область \bar{Q} может быть записана в виде системы отрезков

$$\bar{Q} = \bigcup_{j=1}^q [a^j, b^j], \quad (5.31)$$

на каждом из которых функция неположительна. Для примера рассмотрим рис. 5.3, который отражает возможные случаи поведения непрерывной функции, порождающие области неположительности в виде отрезков (помечены серым цветом), включая касание оси x в точке. Самый правый отрезок, отмеченный пунктиром, соответствует ситуации, когда функция на данном отрезке равна нулю во всех его точках.

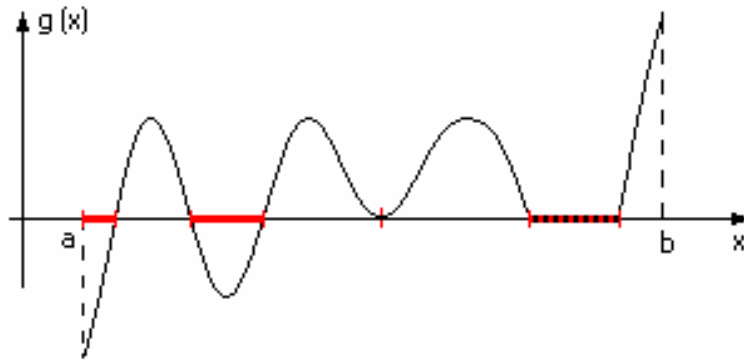


Рис. 5.3. Различные способы порождения областей неположительности непрерывной функции

В системе (5.31) число q отрезков может быть бесконечным. В качестве примера подобной ситуации приведем функцию

$$g(x) = \begin{cases} x \sin \frac{1}{x}, & x > 0, \\ 0, & x = 0, \end{cases}$$

рассматриваемую на отрезке $[0, 1]$.

Таким образом, в случае непрерывной функции $G(y)$ проекция из (5.26) есть множество вида (5.31), т.е.

$$\Pi_i(u_{i-1}) = \bigcup_{j=1}^{q_i} [a_i^j, b_i^j], \quad (5.32)$$

где количество отрезков q_i и их границы a_i^j, b_i^j , $1 \leq j \leq q_i$, зависят от вектора u_{i-1} , т.е.

$$q_i = q_i(u_{i-1}), \quad a_i^j = a_i^j(u_{i-1}), \quad b_i^j = b_i^j(u_{i-1}). \quad (5.33)$$

Если область Y такова, что для всех $1 \leq i \leq N$ удастся указать явные (аналитические) выражения для величин q_i, a_i^j, b_i^j как функций $u_{i-1} \in Y_{i-1}$, тогда область Y называется *областью с вычислимой границей*. Образцы таких областей приведены в примерах 5.1 и 5.2. Для построения данных областей необходимо уметь аналитически находить все корни функций $G^i(u_{i-1}, y_i)$ по соответствующим координатам y_i .

В общем случае, однако, нахождение всех корней непрерывной функции является сложной задачей, не разрешаемой аналитически, и в этом случае можно попытаться найти искомые корни численно. Рассмотрим, к примеру, типовую задачу (5.30) Для отыскания всех корней непрерывной функции $g(x)$ можно предложить численное решение эквивалентной задачи оптимизации

$$|g(x)| \rightarrow \min, x \in [a, b], \quad (5.34)$$

в которой корни $g(x)$ являются точками глобального минимума. Для решения этой задачи могли бы быть использованы характеристические алгоритмы глобального поиска, обеспечивающие сходимость ко всем глобально-минимальным точкам.

Другой подход к учету ограничений в задачах оптимизации (индексный метод), не требующий решения вспомогательных задач (5.34), рассмотрен в параграфе 2.5 предыдущей главы.

Практически важным частным случаем задачи (5.1) – (5.3) является случай $Y = D$, когда функциональные ограничения (5.2) отсутствуют. В данной ситуации $G(y) \equiv 0$ в области D , а из (5.11) следует, что функции $G^i(u_i) \equiv 0$, $u_i \in D_i$. Тогда согласно (5.17)

$$\Pi_i(u_{i-1}) = [a_i, b_i], \quad (5.35)$$

где a_i, b_i - константы.

Другим важным частным случаем является случай *выпуклых* ограничений.

Определение 5.1. Функция $g(y)$ называется *выпуклой (вниз)* в выпуклой области Y , если для любых $y', y'' \in Y$ и для любых $\alpha \in [0,1]$ выполняется

$$g(\alpha y' + (1 - \alpha)y'') \leq \alpha g(y') + (1 - \alpha)g(y'') \quad (5.36)$$

Заметим, что это определение повторяет определение 1.9 из первой главы.

Теорема 5.1. Если в задаче (5.1) – (5.3) ограничения $g_j(y)$, $1 \leq j \leq m$, выпуклы, функция $G(y)$ выбирается согласно (5.6) или (5.7), тогда проекции $\Pi_{i+1}(u_i)$ из (5.10) либо пусты, либо имеют вид

$$\Pi_{i+1}(u_i) = [a_{i+1}^1(u_i), b_{i+1}^1(u_i)], \quad (5.37)$$

т.е. область одномерного поиска в (5.26) состоит из одного отрезка.

ДОКАЗАТЕЛЬСТВО.

1). Покажем вначале, что область Y является выпуклой. Возьмем произвольные $y', y'' \in Y$ и рассмотрим отрезок $\alpha y' + (1 - \alpha)y''$, $\alpha \in [0,1]$. Очевидно, что в силу выпуклости гиперпараллелепипеда D данный отрезок целиком в нем содержится. Кроме того, для любого $\alpha \in [0,1]$ и любого j , $1 \leq j \leq m$, в силу выпуклости $g_j(y)$

$$g_j(\alpha y' + (1 - \alpha)y'') \leq \alpha g_j(y') + (1 - \alpha)g_j(y'') \leq 0,$$

поскольку $g_j(y') \leq 0$ и $g_j(y'') \leq 0$ из-за допустимости точек y' и y'' . Следовательно, $\alpha y' + (1 - \alpha)y'' \in Y$ при любом $\alpha \in [0,1]$, что и доказывает выпуклость Y .

2). Теперь убедимся, что проекции Y_i из (5.15) также выпуклы.

Возьмем два произвольных вектора $u'_i, u''_i \in Y_i$. Тогда существуют вектора v'_i, v''_i такие, что $y' = (u'_i, v'_i), y'' = (u''_i, v''_i) \in Y$. Но тогда для любого $\alpha \in [0,1]$ вектор $\alpha y' + (1-\alpha)y'' = (\alpha u'_i + (1-\alpha)u''_i, \alpha v'_i + (1-\alpha)v''_i) \in Y$, т.е. вектор $\alpha u'_i + (1-\alpha)u''_i \in Y_i$.

3). Покажем, что функция (5.6) выпукла в D .

Для любой точки $\alpha y + (1-\alpha)z$, $\alpha \in [0,1]$, $y, z \in D$, существует номер s такой, что

$$\begin{aligned} G(\alpha y + (1-\alpha)z) &= g_s(\alpha y + (1-\alpha)z) \leq \alpha g_s(y) + (1-\alpha)g_s(z) \leq \\ &\leq \alpha \max\{g_1(y), \dots, g_m(y)\} + (1-\alpha) \max\{g_1(z), \dots, g_m(z)\} = \\ &= \alpha G(y) + (1-\alpha)G(z) \end{aligned}$$

4). Докажем выпуклость функций $G^i(u_i)$ из (5.11) в областях D_i .

Прежде всего, заметим, что функция $G^N(y) \equiv G(y)$ выпукла в $D_N = D$. Предположим теперь выпуклость функции $G^{i+1}(u_{i+1})$ в D_{i+1} , $1 \leq i \leq N-1$. Возьмем произвольные $u'_i, u''_i \in D_i$. Существуют $y'_{i+1}, y''_{i+1} \in [a_{i+1}, b_{i+1}]$ такие, что

$$\begin{aligned} G^i(u'_i) &= G^{i+1}(u'_{i+1}), \quad u'_{i+1} = (u'_i, y'_{i+1}), \\ G^i(u''_i) &= G^{i+1}(u''_{i+1}), \quad u''_{i+1} = (u''_i, y''_{i+1}), \end{aligned}$$

Выберем произвольное $\alpha \in [0,1]$ и обозначим $y^*_{i+1} = \alpha y'_{i+1} + (1-\alpha)y''_{i+1}$ и $u^\alpha_i = \alpha u'_i + (1-\alpha)u''_i$. Тогда существует $y^\alpha_{i+1} \in [a_{i+1}, b_{i+1}]$ такой, что

$$\begin{aligned} G^i(u^\alpha_i) &= G^{i+1}(u^\alpha_i, y^\alpha_{i+1}) \leq G^{i+1}(u^\alpha_i, y^*_{i+1}) \leq \alpha G^{i+1}(u'_i, y'_{i+1}) + (1-\alpha)G^{i+1}(u''_i, y''_{i+1}) = \\ &= \alpha G^i(u'_i) + (1-\alpha)G^i(u''_i) \end{aligned}$$

5). Так как функция $G^{i+1}(u_{i+1})$ выпукла по совокупности переменных u_{i+1} , то она выпукла и по переменной y_{i+1} при фиксированном u_i , т.е. является одномерной выпуклой функцией аргумента y_{i+1} .

Если при этом $G^{i+1}(u_i, y_{i+1}) > 0$ для всех $y_{i+1} \in [a_{i+1}, b_{i+1}]$, тогда $\Pi_{i+1}(u_i) = \emptyset$.

Пусть существует y_{i+1} такой, что $G^{i+1}(u_i, y_{i+1}) \leq 0$. Тогда множество неположительности функции $G^{i+1}(u_i, y_{i+1})$ совпадает с множеством точек глобального минимума функции $\tilde{G}^{i+1}(u_{i+1}) = \max\{0; G^{i+1}(u_i, y_{i+1})\}$, причем функция $\tilde{G}^{i+1}(u_{i+1})$ выпукла и непрерывна как функция максимума от двух выпуклых непрерывных функций. Из выпуклого анализа известно [5], что множество оптимальных точек выпуклой функции на выпуклом множестве также выпукло. В одномерном случае выпуклое множество – это отрезок или (полу)интервал. Так как функция $\tilde{G}^{i+1}(u_{i+1})$ непрерывна по y_{i+1} , то множество ее оптимальных точек на $[a_{i+1}, b_{i+1}]$ является отрезком.

Теорема доказана.

Проекция $\Pi_{i+1}(u_i)$ могут являться отрезками (5.37) не только в случае выпуклых ограничений, но и в более общем случае монотонно унимодальных ограничений [22], которые могут породить невыпуклые множества Y .

Определение 5.2. Пусть $h(\gamma) = (h_1(\gamma), \dots, h_N(\gamma)), \gamma \in [0,1]$, есть параметризованная кривая, соединяющая две точки гиперпараллелепипеда D . Эта кривая называется монотонной, если каждая ее координатная функция $h_i(\gamma)$ является монотонной (одномерной) функцией аргумента $\gamma \in [0,1]$.

Определение 5.3. Функция $g(y)$ называется монотонно унимодальной в D , если для любых двух точек из D существует монотонная кривая h , соединяющая эти точки и такая, что неотрицательная функция

$$g_+(\gamma) = \max\{0; g(h(\gamma))\}, \lambda \in [0,1],$$

либо является унимодальной (вниз), либо она тождественно равна нулю в подынтервале $[\gamma_1, \gamma_2] \subset [0,1]$ и строго монотонна на отрезках $[0, \gamma_1]$ и $[\gamma_2, 1]$, убывая на первом из них и возрастаая на втором.

Определение 5.4. Ограничения (5.2) называются монотонно унимодальными в D в совокупности, если для любой пары точек из D существует единая для всех функций $g_j(y)$ кривая h , обеспечивающая монотонную унимодальность каждого из ограничений.

Примеры.

1. Любая унимодальная (вниз) функция является монотонно унимодальной.
2. Любая выпуклая функция является монотонно унимодальной в D . В качестве кривой $h(\gamma)$ можно взять отрезок. Более того, любой набор выпуклых на D функций будет монотонно унимодальным в совокупности.
3. Пусть

$$D = \{y \in R^2 : y_i \in [0,1], i = 1,2\}$$

Рассмотрим два ограничения: $g_1(y) = y_1^2 + y_2^2 - 0.81$ и $g_2(y) = 0.25 - y_1^2 - y_2^2$. Данные ограничения, являясь монотонно унимодальными в совокупности, порождают в (5.2) невыпуклую допустимую область Y .

Теорема 5.2. Пусть ограничения $g_j(y)$ из (5.2) монотонно унимодальны в совокупности в области D и функция $G(y)$ строится согласно (5.6) или (5.7). Тогда

- 1) $G(y)$ монотонно унимодальна в D ;
- 2) функции $G^i(u_i)$, $1 \leq i \leq N$, монотонно унимодальны в D_i ;
- 3) любая проекция (5.10) имеет вид (5.37).

Доказательство теоремы может быть найдено в [22].

5.4.2. Свойства целевых функций в одномерных подзадачах

Целевой функцией в подзадаче (5.26) является функция $f^i(u_{i-1}, y_i)$ при фиксированном u_{i-1} , и для решения подзадач (5.26) определяющим является характер зависимости функции f^i от переменной y_i .

Рассмотрим класс задач, в которых функция $f(y)$ является сепарабельной, т.е.

$$f(y) = \sum_{i=1}^N f_i(y_i), (5.38)$$

а функциональные ограничения $g_j(y)$ отсутствуют, т.е. $Y = D$.

Тогда, как следует из (5.19), (5.38),

$$\min_{y \in Y} f(y) = \sum_{i=1}^N \min_{y_i \in [a_i, b_i]} f_i(y_i),$$

т.е. для решения многомерной задачи требуется решение N независимых одномерных подзадач. Для этого класса задач порядок роста сложности с ростом размерности линейный.

Предположим теперь, что функция $f(y)$ удовлетворяет в области Y условию Липшица с константой $L > 0$, т.е. для любых $y', y'' \in Y$

$$|f(y') - f(y'')| \leq L \|y' - y''\|. (5.39)$$

Естественным в этом случае является вопрос: а будут ли удовлетворять условию Липшица функции $f^i(u_i)$. Оказывается, ответ не всегда положительный.

Рассмотрим следующий пример. Пусть в задаче (5.1) – (5.3) целевая функция $f(y)$ удовлетворяет условию Липшица (5.39), а область Y имеет вид

$$Y = \{y \in R^2 : y_1^2 + y_2^2 - 1 \leq 0\}$$

Тогда функция $f^2(y) \equiv f(y)$, естественно, является липшицевой с константой L по координате y_2 . Однако функция $f^1(y_1)$ условию (5.38) уже не подчиняется. В [7] показано, что эта функция удовлетворяет обобщенному условию Липшица (условию Гельдера) в метрике $\rho(y'_1, y''_1) = \sqrt{|y'_1 - y''_1|}$ с константой $L_1 = L(1 + \sqrt{2})$, т.е. условию

$$|f(y'_1) - f(y''_1)| \leq L_1 \sqrt{|y'_1 - y''_1|}$$

Следующая теорема [7] устанавливает достаточные условия липшицевости функций $f^i(u_i)$.

Теорема 5.3. Пусть функция $f(y)$ является липшицевой с константой L в выпуклой области Y из (5.2) и граничные пары (5.37) являются кусочно-линейными функциями вида

$$a_{i+1}(u_i) = \max_{1 \leq v \leq p_i} \{\alpha_i^v u_i + A_i^v\}, (5.40)$$

$$b_{i+1}(u_i) = \max_{1 \leq v \leq q_i} \{\beta_i^v u_i + B_i^v\}, (5.41)$$

где $\alpha_i^v u_i, \beta_i^v u_i$ - есть скалярные произведения векторов из R^i и A_i^v, B_i^v - константы. Тогда функции $f^i(u_i), u_i \in Y_i$, являются липшицевыми с константами $L_i, 1 \leq i \leq N$, где

$$L_N = L, \quad L_i = L \prod_{j=i}^{N-1} (1 + \lambda_j), \quad 1 \leq i \leq N-1, \quad \lambda_j = \max \left\{ \max_{1 \leq v \leq p_i} \|\alpha_j^v\|, \max_{1 \leq v \leq q_i} \|\beta_j^v\| \right\}$$

ДОКАЗАТЕЛЬСТВО теоремы приведено в [7].

Комментарии к теореме.

1. Представление граничных пар в виде (5.40), (5.41) имеет место, если область является выпуклым многогранником.

2. Если область Y является гиперпараллелепипедом, то все вектора α_i^v, β_i^v нулевые, и поэтому $L_i = L, 1 \leq i \leq N$.

Основной вывод, который следует из обсуждения липшицевости, состоит в том, что на свойства целевых функций одномерных подзадач существенное влияние оказывают не только свойства исходной целевой функции $f(y)$, но и вид допустимой области Y .

В самом простом случае, когда множество Y является гиперпараллелепипедом, все функции $f^i(u_i)$ будут удовлетворять условию Липшица с той же константой, что и функция $f(y)$.

Если ограничения $g_j(y)$ являются кусочно-линейными выпуклыми функциями, то в этом случае липшицевость функций $f^i(u_i)$ сохраняется, но константа Липшица для этих функций, вообще говоря, увеличивается, что ухудшает оптимизационные свойства одномерных подзадач.

Наконец, случай нелинейных ограничений может вообще привести к потере липшицевости.

5.5. Редукция размерности на основе кривых Пеано

Рассмотрим кривую Пеано $y(x), x \in [0, 1]$, отображающую однозначно и непрерывно отрезок вещественной оси на многомерный гиперпараллелепипед (5.3), т.е.

$$D = \{y(x) : 0 \leq x \leq 1\} \quad (5.42)$$

Непрерывность целевой функции $f(y)$ многомерной задачи (5.1)-(5.3) (для простоты – в отсутствие функциональных ограничений (5.2), т.е. $Y = D$) обеспечивает соотношение

$$\min_{y \in D} f(y) = \min_{x \in [0, 1]} f(y(x)) \quad (5.43)$$

и возможность вместо многомерной задачи (5.1) решать одномерную задачу

$$f(y(x)) \rightarrow \min, x \in [0, 1] \quad (5.44)$$

Вместе с тем эта возможность может быть реализована только при наличии алгоритмического способа реализации развертки $y(x)$. Рассмотрим одну из схем построения кривых Пеано, предложенную Гильбертом.

Указанное построение формирует непрерывное и однозначное отображение $y(x)$ отрезка $[0, 1]$ на стандартный гиперкуб

$$D = \{y \in R^N : -\frac{1}{2} \leq y_i \leq \frac{1}{2}, 1 \leq i \leq N\} \quad (5.45)$$

Рассмотрение стандартного гиперкуба не умаляет общности, поскольку гиперпараллелепипед (5.3) очевидным образом приводится к виду (5.45) линейным преобразованием координат.

Предлагаемое построение приведено, например, в [7] и состоит в следующем.

1. Сначала гиперкуб D из (5.45), длина ребра которого равна 1, разделяется координатными плоскостями на 2^N гиперкубов так называемого *первого разбиения* (с длиной ребра, равной $1/2$), которые занумеруем числами z_1 от 0 до $2^N - 1$, причем гиперкуб первого разбиения с номером z_1 условимся обозначать через $D(z_1)$.

Далее каждый гиперкуб первого разбиения в свою очередь также разбивается на 2^N гиперкубов *второго разбиения* (с длиной ребра, равной $1/4$) гиперплоскостями, параллельными координатным и проходящим через серединные точки ребер гиперкуба, ортогональных к этим гиперплоскостям. При этом гиперкубы второго разбиения, входящие в гиперкуб $D(z_1)$, нумеруются числами z_2 от 0 до $2^N - 1$, причем гиперкуб второго разбиения с номером z_2 , входящий в $D(z_1)$, условимся обозначать через $D(z_1, z_2)$.

Продолжая указанный процесс, можно построить гиперкубы любого p -го разбиения с длиной ребра, равной $(1/2)^p$, которые условимся обозначать через $D(z_1, \dots, z_p)$, причем

$$D(z_1) \supset D(z_1, z_2) \supset \dots \supset D(z_1, \dots, z_p)$$

и $0 \leq z_j \leq 2^N - 1, 1 \leq j \leq p$.

2. Теперь осуществим деление отрезка $[0, 1]$ на 2^N равных частей, каждую из которых в свою очередь также разделим на 2^N равных частей и т.д., причем элементы каждого разбиения будем нумеровать слева направо числами z_j (j -номер разбиения) от 0 до $2^N - 1$. При этом интервалы p -го разбиения условимся обозначать как $d(z_1, \dots, z_p)$, где, например, $d(z_1, z_2)$ обозначает интервал второго разбиения с номером z_2 , являющийся частью интервала $d(z_1)$ первого разбиения с номером z_1 . Заметим, что

$$d(z_1) \supset d(z_1, z_2) \supset \dots \supset d(z_1, \dots, z_p)$$

и длина интервала $d(z_1, \dots, z_p)$ равна $(1/2)^{pN}$.

Предполагается, что интервал $d(z_1, \dots, z_p)$ содержит свой левый конец; он содержит правый конец тогда и только тогда, когда

$$z_1 = z_2 = \dots = z_p = 2^N - 1$$

3. Примем, что точка $y(x) \in D$, соответствующая точке $x \in [0, 1]$, при любом $p \geq 1$ содержится в гиперкубе $D(z_1, \dots, z_p)$, если x принадлежит интервалу $d(z_1, \dots, z_p)$, т.е.

$$x \in d(z_1, \dots, z_p) \rightarrow y(x) \in D(z_1, \dots, z_p)$$

Построенное соответствие $y(x)$ является однозначным.

4. Для любого $x \in d(z_1, \dots, z_p)$ справедливо, что

$$x - \left(\frac{1}{2}\right)^{pN} \leq \sum_{j=1}^p z_j \left(\frac{1}{2}\right)^{jN} \leq x + \left(\frac{1}{2}\right)^{pN}$$

Следовательно, если представить x в виде двоичного числа с фиксированной запятой, т.е.

$$x = \sum_{i=1}^{\infty} \alpha_i \left(\frac{1}{2}\right)^i, \quad (5.46)$$

где α_i есть двоичные цифры ($\alpha_i = 0$ или $\alpha_i = 1$), то по первым pN цифрам $\alpha_1, \dots, \alpha_{pN}$ этого числа можно указать гиперкуб p -го разбиения $D(z_1, \dots, z_p)$, содержащий точку $y(x)$, поскольку

$$z_j = \sum_{i=0}^{N-1} \alpha_{Nj-i} (2^i), \quad 1 \leq j \leq p, \quad (5.47)$$

где α_{Nj-i} из (5.46). Таким образом, точка $y(x)$ может быть оценена с точностью $(1/2)^{p+1}$ по каждой координате, если известно (5.47).

5. Чтобы построенное отображение было вдобавок непрерывным, наложим требования на порядок нумерации гиперкубов каждого разбиения.

Определим вектор $z^p = (z^1, \dots, z^p)$, соответствующий гиперкубу $D(z_1, \dots, z_p)$ p -го разбиения, и условимся, что вектор z^p предшествует вектору \bar{z}^p , если при $p \geq 1$ либо $z_1 < \bar{z}_1$, либо существует такое k , $1 \leq k < p$, что $z_j = \bar{z}_j, 1 \leq j \leq k$, и $z_{k+1} < \bar{z}_{k+1}$.

Введенное отношение устанавливает совершенный строгий порядок на множестве различных векторов z^p . Минимальным элементом при этом является вектор $(0, \dots, 0)$, а максимальным - вектор $(2^N - 1, \dots, 2^N - 1)$. Порядок данного вида фактически задает лексикографическое упорядочение векторов z^p .

Будем говорить, что векторы z^p и \bar{z}^p (и соответственно гиперкубы $D(z_1, \dots, z_p)$ и $D(\bar{z}_1, \dots, \bar{z}_p)$) являются смежными, если один из них предшествует другому и не существует вектора \tilde{z}^p такого, что

$$z^p \prec \tilde{z}^p \prec \bar{z}^p \quad \text{или} \quad \bar{z}^p \prec \tilde{z}^p \prec z^p$$

где \prec - символ отношения предшествования.

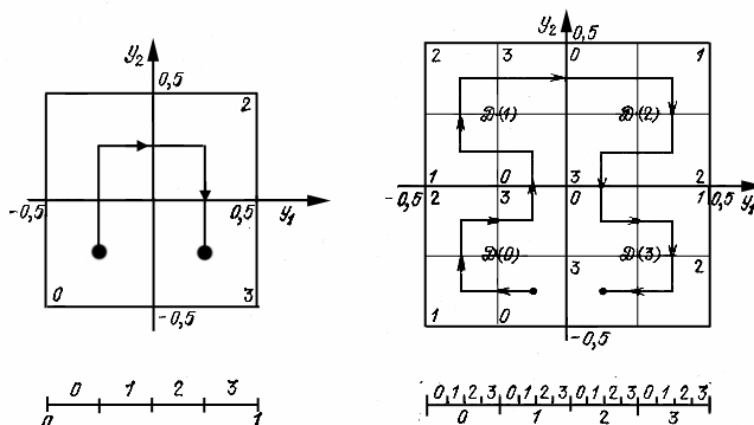


Рис. 5.4. Нумерация гиперкубов при построении первых двух приближений кривой Пеано

Для непрерывности построенного отображения достаточно выбрать такую нумерацию, чтобы смежные гиперкубы любого p -го разбиения ($p \geq 1$) имели общую грань. Тогда, устремляя p к бесконечности, в пределе получим непрерывную кривую Пеано, заполняющую гиперкуб D .

Существуют различные способы требуемой нумерации, детальное описание которых можно найти, например, в [7]. Продемонстрируем на примере, взятом из [11], один из вариантов такой нумерации (рис. 5.4, 5.5).

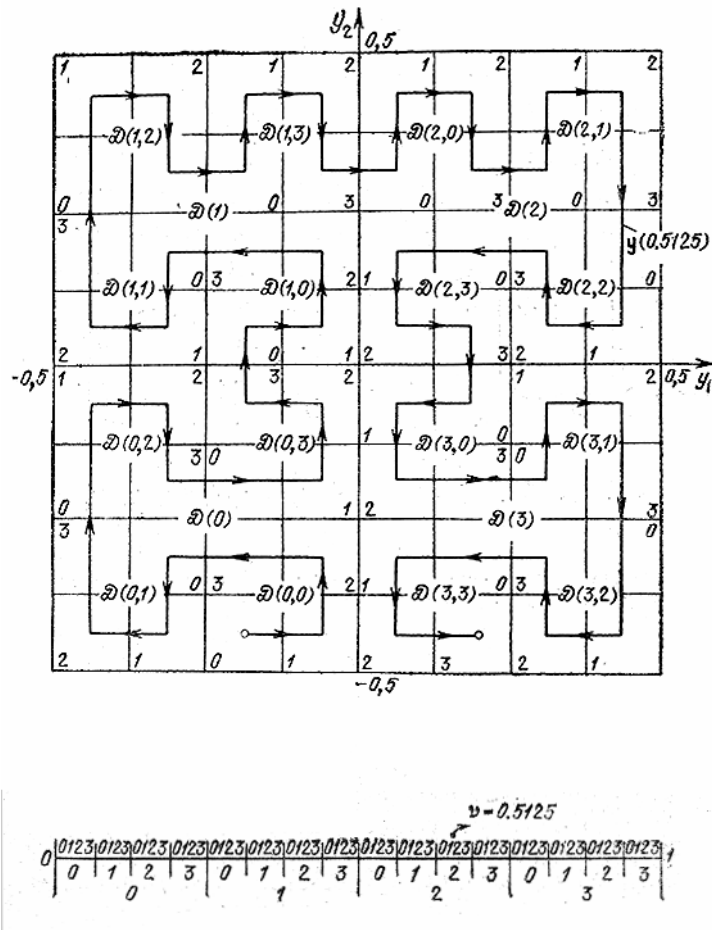


Рис. 5.5. Следующее приближение кривой Пеано

Построенное отображение Пеано не обеспечивает взаимно однозначное соответствие между отрезком $[0,1]$ и N -мерным гиперкубом (5.45), потому что точка $y \in D = \{y(x) : x \in [0,1]\}$ может иметь несколько прообразов в отрезке $[0,1]$, но не более, чем 2^N . Кратность точек $y \in D$ при соответствии $y(x)$ является фундаментальным свойством, отражающим сущность понятия размерности: отрезок $[0,1]$ и гиперкуб D являются равномошными множествами, первое из них можно однозначно отобразить на второе, но если такое отображение непрерывно, то оно не может быть взаимно однозначным. При этом размерность N гиперкуба определяет возможную кратность (2^N) для отображения $y(x)$. Если вернуться к использованию разверток для целей оптимизации, то это означает, что точка y^* глобального минимума функции $f(y)$ в гиперкубе D может порождать несколько точек x_s^* глобального минимума функции $f(y(x))$ в отрезке $[0,1]$ таких, что $y^* = y(x_s^*)$.

Теорема 5.4. [7] Пусть $f(y)$ есть липшицева с константой L функция, определенная в области D из (5.45). Тогда функция $f(y(x)), x \in [0,1]$, где $y(x)$ есть непрерывное однозначное отображение со свойствами 1-5, удовлетворяет условию

$$|f(y(x')) - f(y(x''))| \leq L_0 \sqrt[N]{|x' - x''|} \quad (5.48)$$

где $x', x'' \in [0,1]$ и $L_0 = 4L\sqrt{N}$.

ДОКАЗАТЕЛЬСТВО. Пусть $|x' - x''| > 0$. Тогда существует такое целое число $p \geq 0$, что

$$(1/2)^{(p+1)N} \leq |x' - x''| \leq (1/2)^{pN} \quad (5.49)$$

При этом точки x' и x'' будут принадлежать либо одному и тому же интервалу p -го разбиения $d(z_1, \dots, z_p)$, либо двум разным, но смежным интервалам. Следовательно, их образы $y(x')$ и $y(x'')$ принадлежат либо одному и тому же соответствующему гиперкубу p -го разбиения, либо двум разным гиперкубам p -го разбиения, имеющим общую грань. Тогда

$$\|y(x') - y(x'')\| \leq \sqrt{N} (1/2)^{p-1},$$

что вследствие (5.49) влечет

$$\|y(x') - y(x'')\| \leq 4\sqrt{N} \sqrt[N]{|x' - x''|}, \quad (5.50)$$

откуда следует (5.48), поскольку функция $f(y)$ является липшицевой с константой L .

Теорема доказана.

Таким образом, функция $f(y(x)), x \in [0,1]$, удовлетворяет обобщенному условию Липшица (5.48) с константой L_0 в метрике

$$\rho(x', x'') = \sqrt[N]{|x' - x''|} \quad (5.51)$$

Тогда для решения одномерной задачи (5.44) можно применить характеристические алгоритмы, ориентированные на минимизацию липшицевых функций, заменив в их решающих правилах "евклидово" расстояние $|x' - x''|$ на расстояние (5.51). Так, для информационно-статистического алгоритма глобального поиска АГП его модификация – многомерный обобщенный алгоритм глобального поиска (МОАГП) – содержит характеристику

$$R(i) = m \sqrt[N]{x_i - x_{i-1}} + \frac{(z_i - z_{i-1})^2}{m \sqrt[N]{x_i - x_{i-1}}} - 2(z_i + z_{i-1}), \quad (5.52)$$

где m из (4.44), $M = \max_{1 \leq i \leq \tau} \frac{|z_i - z_{i-1}|}{\sqrt[N]{x_i - x_{i-1}}}$, а $z_j = f(y(x_j)), 0 \leq j \leq \tau$.

Точка очередного испытания в МОАГП вычисляется согласно выражению

$$x^{k+1} = \frac{(x_t + x_{t-1})}{2} - \frac{\text{sign}(z_t - z_{t-1})}{2r} \left\{ \frac{r|z_t - z_{t-1}|}{m} \right\}^N. \quad (5.53)$$

Нетрудно показать, что МОАГП удовлетворяет условиям теорем 4.1, 4.4, а при $m > 16L\sqrt{N}$ и достаточным условиям сходимости теоремы 4.5.

Напомним, что кривая Пеано, построенная по описанной схеме, задается через предельный переход, поэтому в практической ситуации реализуемо лишь некоторое приближение к этой кривой. Рассмотрим несколько способов построения таких приближений и их свойства.

Пеаноподобная развертка. Сначала остановимся на кусочно-линейной развертке из [7]. Согласно выполненной ранее схеме построения 2^{pN} центров $y(z_1, \dots, z_p)$ гиперкубов p -го разбиения $D(z_1, \dots, z_p)$ образуют равномерную ортогональную сетку $H(p, N)$ в гиперкубе D из (5.45). Шаг этой сетки (по любой координате) равен 2^{-p} . Введем следующую нумерацию узлов данной сетки. Занумеруем слева направо индексом i все интервалы, составляющие p -е разбиение отрезка $[0, 1]$, т.е.

$$d(z_1, \dots, z_p) = [x_i, x_{i+1}), \quad 0 \leq i \leq 2^{pN} - 1,$$

где через x_i обозначен левый конец интервала, имеющего номер i . Будем считать, что центр гиперкуба имеет тот же номер i , что и соответствующий этому гиперкубу (согласно введенной исходной связи между номерами гиперкубов и интервалов) интервал $d(z_1, \dots, z_p)$, т.е.

$$y_i = y(z_1, \dots, z_p), \quad 0 \leq i \leq 2^{pN} - 1.$$

При такой нумерации центры y_i и y_{i+1} необходимо соответствуют смежным гиперкубам, имеющим общую грань.

Рассмотрим отображение $q(x)$ отрезка $[0, 1]$ в гиперкуб (5.45), определяемое соотношениями

$$q(x) = y_i + (y_{i+1} - y_i) \left(\frac{w(x) - x_i}{x_{i+1} - x_i} \right), \quad (5.54)$$

$$x_i \leq w(x) \leq x_{i+1},$$

$$w(x) = x \left(1 - (1/2)^{pN} \right), \quad 0 \leq x \leq 1. \quad (5.55)$$

(величины $q(x), y_i, y_{i+1}$ суть N -мерные вектора).

Очевидно, что участок кривой $q(x)$, соединяющий узлы y_i и y_{i+1} , является линейным, а это значит, что $q(x)$, $0 \leq x \leq 1$, есть кусочно-линейная кривая, соединяющая узлы y_i , $0 \leq i \leq 2^{pN} - 1$ в порядке их нумерации. Эту кривую в дальнейшем будем называть *пеаноподобной кусочно-линейной разверткой* отрезка в гиперкубе D из (5.45), поскольку она является приближением (с точностью не хуже 2^{-p} по каждой координате к кривой Пеано. На рис. 5.6 изображен образ отрезка $[0, 1]$ при соответствии $q(x)$ для двухмерного пространства и разбиения 3 порядка.

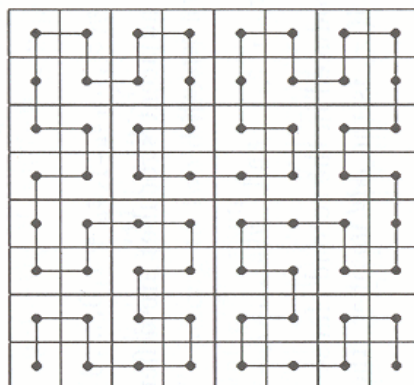


Рис. 5.6. Образ отрезка $[0,1]$ при отображении пeanоподобной развертки при разбиении третьего порядка

Введенное отображение $q(x)$ сопоставляет минимизируемой многомерной функции $f(y)$ одномерную функцию $f(q(x))$, и если функция $f(y)$ удовлетворяет условию Липшица (с константой L), то

$$\min_{y \in D} f(y) \leq \min_{x \in [0,1]} f(q(x)) + (L\sqrt{N}) \left(\frac{1}{2}\right)^{p+1} \quad (5.56)$$

и приближенное решение многомерной задачи сводится к решению задачи одномерной, где функция $f(q(x))$ удовлетворяет обобщенному условию Липшица в метрике (5.51) (доказательство аналогично доказательству Теоремы 5.4). Из этого следует, что для решения одномерной задачи можно использовать обобщенные на метрику (5.51) алгоритмы типа алгоритма МОАГП (5.52), (5.53).

Простые кусочно-линейные развертки. Наряду с пeanоподобной кусочно-линейной разверткой можно рассмотреть более простые непрерывные развертки, покрывающие сетку $H(p, N)$, для которых остается в силе соотношение (5.56). В качестве таких разверток можно взять, например, так называемую "телевизионную" кривую или спиральную развертку (см. рис. 5.7).

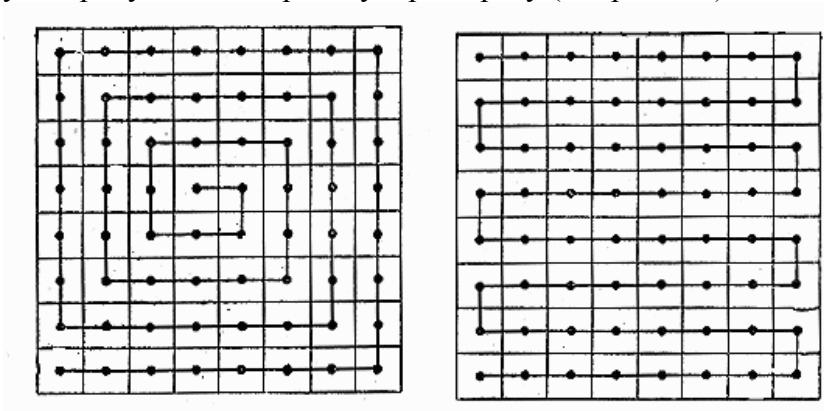


Рис. 5.7. Примеры простых кусочно-линейных разверток

Вместе с тем простые кусочно-линейные развертки $t(x)$ имеют один существенный недостаток, обусловленный тем обстоятельством, что соответствующая кусочно-линейная кривая содержит линейные участки, покрывающие большое число (до 2^p) узлов сетки $H(p, N)$. Это приводит к тому,

что даже для метрики (5.51) не существует единой для всех значений $p > 1$ константы, обеспечивающей выполнение обобщенного условия Липшица для функции $f(t(x))$, даже если функция $f(y)$ является липшицевой. Более детально, при конкретном p функция $f(t(x))$ удовлетворяет условию (5.48) с некоторой константой L_p , но при $p \rightarrow \infty$ константа L_p становится неограниченно большой, что приводит к увеличению плотности испытаний метода оптимизации и ухудшению показателей эффективности по сравнению с пеаноподобной разверткой.

Неинъективная развертка. Пеаноподобная кусочно-линейная развертка $q(x)$, превосходя простые развертки, тем не менее, не обладает двумя важными свойствами, которые есть у кривой Пеано, ею приближаемой.

Первый недостаток связан с тем, что сетка $H(p + \lambda, N)$, имеющая шаг $(1/2)^{p+\lambda}$, не содержит узлов более грубой сетки. Поэтому при уточнении сетки в процессе поиска невозможно учесть уже вычисленные на сетке $H(p, N)$ значения функции, и оптимизацию надо начинать сначала.

Второй недостаток развертки $q(x)$ связан с тем обстоятельством, что она является *взаимно однозначным* соответствием между отрезком $[0, 1]$ и множеством $\{q(x) : 0 \leq x \leq 1\}$, лежащим в области D , хотя кривая Пеано таким свойством не обладает. Напомним, что кривая Пеано, являясь однозначным отображением отрезка $[0, 1]$ в N -мерный гиперкуб D из (5.45), вместе с тем взаимно однозначным соответствием не является, потому что точка $y \in D = \{y(x) : x \in [0, 1]\}$ может иметь несколько прообразов в отрезке $[0, 1]$, но не более, чем 2^N экземпляров. Это означает, что близкие точки x', x'' из отрезка $[0, 1]$ переходят в близкие точки y', y'' в гиперкубе D , но две близкие точки y', y'' из гиперкуба D могут иметь существенно "далекие" прообразы $x', x'' \in [0, 1]$.

Применительно к глобальной оптимизации это означает следующее. Любая точка глобального минимума $y^* \in D$ может иметь несколько прообразов $x_i^* \in [0, 1], 1 \leq i \leq \nu \leq 2^N$, и некоторые из этих прообразов располагаются далеко друг от друга. Следовательно, результаты поиска на последовательности $\{y^k\} = \{y(x^k)\}$, порождаемой оптимизационной схемой с редукцией размерности по кривым Пеано, для точек испытаний y^n из некоторой сколь угодно малой окрестности y^* (близких в смысле принадлежности к этой окрестности) представляются соответствующими парами (x_i, z_i) в поисковой информации ω_k , но находятся *только в некоторых* окрестностях прообразов $x_i^*, 1 \leq i \leq \nu$. Таким образом, имеются прообразы x_i^* с окрестностями, в которых потеряна информация об уже проведенных испытаниях в том смысле, что в этих окрестностях есть точки, соответствующие точкам y^n , в которых получены значения $z^n = f(y^n)$, но эти значения не учтены в поисковой информации ω_k . Из-за этого в таких окрестностях приходится проводить дополнительные испытания, так как все точки $x_i^*, 1 \leq i \leq \nu$ являются предельными точками последовательности $\{x^k\} \subset [0, 1]$, генерируемой алгоритмом оптимизации при решении задачи (5.44).

В связи с этим можно построить другую развертку, которая отображает некоторую равномерную сетку в отрезке $[0, 1]$ на сетку $U(p, N)$, составленную из

вершин гиперкубов $D(z_1, \dots, z_p)$ p -го разбиения гиперкуба D . Шаг этой сетки по любой координате равен $(1/2)^p$, общее число различных узлов равно $(2^p + 1)^N$ и поскольку вершины гиперкубов p -го разбиения являются также вершинами гиперкубов любого следующего разбиения $p + \lambda$, то справедливо включение

$$U(p, N) \subset U(p + \lambda, N).$$

Вводя соответствующее равномерное разбиение отрезка $[0, 1]$ и устанавливая определенное соответствие между узлами данного разбиения и вершинами гиперкубов, задающими сетку, мы получим новую развертку, которая также будет являться приближением к кривой Пеано, но в отличие от развертки $q(x)$ будет иметь кратные точки (узел сетки $U(p, N)$ может иметь до 2^N прообразов в отрезке $[0, 1]$). Развертку $s(x)$ будем называть *пеаноподобной неинъективной разверткой*. Детальное построение неинъективной развертки приведено в [7].

Применяя неинъективную развертку в глобальной оптимизации, можно теперь вычислять и включать в поисковую информацию все прообразы каждой точки y^k из последовательности $\{y^k\} = \{y(x^k)\}$, которых может быть достаточно много (до 2^N).

Множественная развертка.

Применяя неинъективную развертку, мы пополняем поисковую информацию ω_k значительным объемом данных (до 2^N пар (x_i, z_i)) на каждую точку многомерного испытания, где x_i - один из прообразов точки y^k , а $z_i = f(y^k)$. Р.Г.Стронгин [9] предложил новую схему редукции размерности, которая сохраняет информацию о близости точек в многомерном пространстве и за счет этого существенно уменьшает объем данных, включаемых в поисковую информацию ω_k одномерного поиска. Более конкретно, объем данных, пополняющих поисковую информацию после каждой итерации поиска *не зависит от размерности* исходной проблемы, а определяется требуемой точностью решения. Предложенная схема редукции получила название "*множественная развертка*". Кратко опишем основные идеи схемы.

Рассмотрим интервал $[0, L + 1]$ на оси x и семейство гиперкубов

$$D_l = \{y \in R^N : -2^{-1} - 2^{-l} \leq y_j \leq 3 \cdot 2^{-1} - 2^{-l}, 1 \leq j \leq N\}, \quad 0 \leq l \leq L, \quad (5.57)$$

где каждый гиперкуб $D_l, 1 \leq l \leq L$, получается смещением гиперкуба D_{l-1} вдоль главной диагонали на расстояние, соответствующее смещению 2^{-l} по каждой координате, причем

$$D_0 = \{y \in R^N : -3 \cdot 2^{-1} \leq y_j \leq 2^{-1}, 1 \leq j \leq N\}, \quad (5.58)$$

Рисунок представляет случай $L = N = 2$.

Преобразование

$$y^0(x) = F(y(x)) = 2y(x) - I \quad (5.59)$$

где

$$I = (1, 1, \dots, 1) \in R^N, \quad (5.60)$$

и $y(x)$ является кривой Пеано, сформированной по схеме Гильберта, обеспечивает соответствие $y^0(x)$, непрерывно отображающее единичный интервал $[0,1]$ оси x на гиперкуб D_0 из (5.58), т.е.

$$D_0 = \{y^0(x) : x \in [0, 1]\} \quad (5.61)$$

В представлении (5.59) любой подкуб $D(p, v) = D(z_1, \dots, z_p)$, где $v = (z_1, \dots, z_p)$, некоторого p -го разбиения куба D из (5.45) имеет образ

$$D_0(p, v) = F(D(p, v)), \quad (5.62)$$

который является подкубом p -го разбиения куба D_0 из (5.58). Вследствие (5.62) оба подкуба $D(p, v)$ и $D_0(p, v)$ соответствуют одному и тому же интервалу $d(p, v) = d(z_1, \dots, z_p)$ p -го разбиения единичного интервала $[0,1]$, имеющего длину 2^{-pN} . Заметим, что длина ребра подкуба $D_0(p, v)$ равна $2^{-(p-1)}$, потому что ребро куба D_0 в два раза больше ребра куба D . Последнее утверждение верно для длин ребер любого куба $D_l, 1 \leq l \leq L$, из (5.57). Следовательно, если две точки $x', x'' \in [0, 1]$ удовлетворяют условию

$$|x' - x''| \leq 2^{-(p+1)N}, \quad (5.63)$$

тогда их образы должны удовлетворять неравенству

$$\max \{|y_j^0(x') - y_j^0(x'')| : 1 \leq j \leq N\} \leq 2^{-(p-1)} \quad (5.64)$$

Далее, введем семейство кривых $y^l(x)$ (называемых развертками), определенных рекуррентными соотношениями

$$y^l(x) = y^{l-1}(x) + 2^{-l}I, \quad 1 \leq l \leq L, \quad (5.65)$$

где I из (5.60). $y^l(x)$ отображает единичный интервал $[0,1]$ на соответствующий гиперкуб $D_l, 1 \leq l \leq L$, из (5.57), т.е.

$$D_l = \{y^l(x) : x \in [0, 1]\} \quad (5.66)$$

Как результат, соответствие

$$Y(x) = y^{\lfloor x \rfloor}(x - \lfloor x \rfloor), \quad x \in [0, L+1], \quad (5.67)$$

где $\lfloor x \rfloor$ обозначает целую часть x , отображает единичный интервал $[l, l+1)$ на гиперкуб $D_l, 0 \leq l \leq L$. Из (5.65), (5.66) следует, что для каждого частичного подкуба $D_0(p, v)$ p -го разбиения гиперкуба D_0 существует подкуб $D_l(p, v)$ p -го разбиения гиперкуба D_l такой, что он может быть получен переносом $D_0(p, v)$ вдоль главной диагонали на расстояние $2^{-1} + 2^{-2} + \dots + 2^{-l}$ (по каждой координате). К примеру, на рисунке подкубы $D_0(2, 5 \cdot 2^{-3})$ и $D_1(2, 5 \cdot 2^{-3})$ заштрихованы и помечены символами D^0 и D^1 .

Как следует из (5.65)-(5.67), если $d(p, v)$ соотносится с $D_0(p, v)$, тогда существует ряд подкубов

$$D_l(p, v) = y^l(d(p, v)), \quad 0 \leq l \leq L,$$

соответствующих подынтервалам

$$d_l(p, v_l) \subset [l, l+1), v = v_l - \lfloor v_l \rfloor, \quad 0 \leq l \leq L,$$

где $d_0(p, v_0) = d(p, v)$, $v = v_0$, и

$$D_l(p, v) = Y(d_l(p, v_l)), \quad 0 \leq l \leq L. \quad (5.68)$$

Заметим, что подкубы $D_l(p, v)$, полученные переносом соответствующих подкубов $D_{l-1}(p, v)$, характеризуются тем же самым параметром v , который является точкой левого конца интервала $d(p, v) = d_0(p, v) \subset [0, 1]$. Так как подкубы $D_l(p, v)$ соотносятся с некоторым интервалом $d_l(p, v_l) \subset [l, l+1)$, точки v_l левых концов этих интервалов связаны с v очевидным соотношением $v = v_l - \lfloor v_l \rfloor$.

Принимая во внимание очевидное включение $D \subset D_l, 0 \leq l \leq L$, можно представить исходный гиперкуб D из (5.45) в любой из следующих форм:

$$D = \{Y(x) : x \in [l, l+1), g_0(Y(x)) \leq 0\}, \quad 0 \leq l \leq L, \quad (5.69)$$

где

$$g_0(y) = \max \{|y_j| - 2^{-1} : 1 \leq j \leq N\} \quad (5.70)$$

и, следовательно, $g_0(y) \leq 0$, если $y \in D$, и $g_0(y) > 0$, если $y \notin D$. Как следует из (5.69), развертка (5.67) покрывает гиперкуб D из (5.45) $L+1$ раз, когда x изменяется от 0 до $L+1$. Именно поэтому данная кривая называется *множественной разверткой*.

Как следствие, любая точка $y \in D$ имеет прообраз x^l в каждом интервале $[l, l+1)$, т.е.

$$y = y^{\lfloor x^l \rfloor}(x - \lfloor x^l \rfloor), \quad x^l \in [l, l+1), \quad 0 \leq l \leq L, \quad (5.71)$$

(см. рис. 5.8). Основные свойства кривой (5.67) определяются следующими утверждениями.

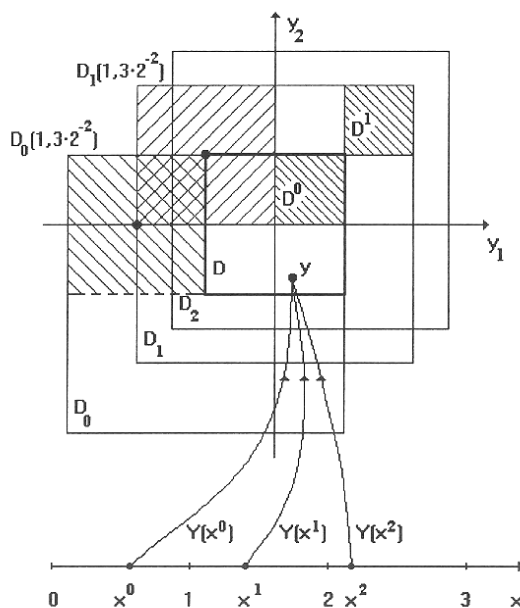


Рис. 5.8. Свойства множественной развертки

Теорема 5.5. Пусть точка y^* из области (5.45) содержится в линейном сегменте с конечными точками $y', y'' \in D$, удовлетворяющими условиям:

$$|y'_j - y''_j| \leq 2^{-p}; \quad y'_i = y''_i = y_i^*, 1 \leq i \leq N, i \neq j, \quad (5.72)$$

где p - некоторое целое число и $1 \leq p \leq L$; т.е. линейный сегмент параллелен j -й оси в R^N . Тогда кривая (5.67) обеспечивает существование целого числа l , $0 \leq l \leq L$, и прообразов $x^*, x', x'' \in [l, l+1)$, удовлетворяющих следующим условиям:

$$y^* = Y(x^*), \quad y' = Y(x'), \quad y'' = Y(x''), \quad (5.73)$$

и

$$\max\{|x' - x^*|, |x'' - x^*|, |x' - x''|\} \leq 2^{-pN}. \quad (5.74)$$

Комментарий к теореме. Условия (5.72) вводят частный тип т.н. 2^{-p} -окрестности точки y^* . Эта окрестность охватывает всех "соседей" заданной точки y^* , отличающихся от нее только j -й координатой, но не более, чем на величину 2^{-p} . Максимально возможное значение p зависит от числа L используемых разверток. Изменяя j в (5.72), можно получить соседей в любом из N координатных направлений.

Как утверждает теорема, два любых соседа из j -й 2^{-p} -окрестности точки y^* имеют по меньшей мере два прообраза в 2^{-pN} -окрестности некоторого прообраза x^* на оси x . Это является способом, которым обратное многозначное отображение $Y^{-1}(y)$ отражает свойство близости в R^N в любых N координатных направлениях.

5.6. Решение задач с ограничениями с использованием разверток

Соединив идею редуцирования размерности посредством разверток со структурой индексного метода учета ограничений (см. пункт 3.4.3 третьей главы), нетрудно получить метод решения общей задачи (5.1)–(5.3) при наличии функциональных ограничений. Для этого достаточно каждой точке $x \in [0, 1]$ приписать индекс $\nu(x) = \nu(y(x))$, где $\nu(y(x))$ — номер первого нарушенного ограничения в точке $y(x)$ (либо $\nu(y(x)) = m + 1$, если все ограничения выполнены), а затем провести решение задачи (5.44), применяя индексный метод (3.54)–(3.59). Иллюстрации для размерности $N=1$ приведены на рис. 3.10, 3.11 в третьей главе.

Детальное изложение алгоритма и особенностей его реализации и применения может быть найдено в [44].

5.7. Компонентные методы

Данный подход к решению многомерных многоэкстремальных задач можно рассматривать как обобщение характеристических принципов построения методов одномерной оптимизации на многомерный случай.

Идея подхода состоит в том, что допустимая область D из (5.2) делится на $\tau = \tau(k) > 0$ частей (компонент) $D_i^k, 1 \leq i \leq \tau(k)$, так, что $D = \bigcup_{i=1}^{\tau(k)} D_i^k$ и их внутренние области попарно не пересекаются. Затем каждой компоненте D_i^k приписывается число $R(i) = R(D_i^k)$, называемое ее характеристикой, и выбирается компонента с наибольшей характеристикой. В этой компоненте проводится очередное испытание (или группа испытаний), а затем происходит ее деление на несколько более "мелких" компонент, так чтобы им принадлежали точки новых проведенных испытаний. При этом число компонент увеличивается, и возникает новое разбиение $D_i^{k+1}, 1 \leq i \leq \tau(k+1)$. Модель поведения функций задачи в пределах каждой компоненты обычно принимается зависящей только от испытаний, проведенных в этой компоненте и параметров метода.

Условия сходимости для каждого подкласса методов компонентного типа обычно устанавливаются авторами методов с учетом специфики данного конкретного подкласса, хотя возможен более общий подход. Здесь следует отметить работу Я.Д. Сергеева, где он обобщил характеристическую теорию сходимости одномерных алгоритмов на компонентную схему реализации многомерных методов, введя понятие "divide-the-best" алгоритмов с определенными требованиями к механизму деления на компоненты. В этой теории установлены [54] основные свойства сходимости, развивающие свойства сходимости характеристического одномерного подхода. В частности, установлено свойство многосторонней (точнее, 2^N – сторонней) сходимости к предельной точке как аналог двусторонней сходимости в одномерном случае. Аналогичный результат по анализу многосторонней сходимости в рамках теории T – представимых методов представлен в пункте 4.5.3 четвертой главы этой книги.

Центральным моментом алгоритмов, исповедующих идею деления области на компоненты, является задание характеристик $R(i)$ компонент D_i^k , учитывающих результаты выполненных в D_i^k испытаний. Новым (и сложным!) моментом, по сравнению с одномерным случаем, является процедура расщепления наилучшей компоненты, а также борьба с возможной избыточностью измерений.

Кратко рассмотрим одни из первых компонентных методов, которые были предложены Я. Пинтером [52], Ю.Г. Евтушенко [26], а также приведем несколько новых схем компонентных методов [19,20,40].

5.7.1. Метод деления на три

Метод был первоначально предложен в работе Ю.Г. Евтушенко и В.А. Ратькина [26] под название *метода половинных делений*, позднее теми же авторами была разработана более экономичная версия, использующая деления на три. В данном разделе рассмотрена именно эта версия. Метод деления на три позволяет решать задачи вида (5.3) как при отсутствии, так и при наличии функциональных ограничений. Предполагается, что компоненты вектор–функции $Q = (f, g_1, \dots, g_m)$ липшицевы с известными константами L_0, L_1, \dots, L_m .

АЛГОРИТМ МЕТОДА ДЕЛЕНИЯ НА ТРИ (принципиальная схема).

ШАГ 0. задается ε — точность решения.

ШАГ 1. Многомерный параллелепипед D разделяется по большему ребру на три равные компоненты в виде гиперпараллелепипедов D_1, D_2, D_3 . В их геометрических центрах y^1, y^2, y^3 выполняются измерения набора функций Q .

ШАГ 2. Для имеющихся компонент D_i по измерениям Q^i в их центрах строятся нижние оценки значений компонент функции Q в следующем виде

$$f^-(D_i) = f^i - L_0 d(D_i)/2, \quad g_j^-(D_i) = g_j^i - L_j d(D_i)/2 \quad (5.75)$$

где $d(D_i)$ — диаметр гиперпараллелепипеда D_i , в принятой метрике. Вычисляются характеристики компонент в виде

$$R(i) = \begin{cases} f^-(D_i), & \max\{g_j^-(D_i) : j=1, \dots, m\} \leq 0 \\ +\infty, & \exists j \in \{1, \dots, m\}, \text{ что } g_j^-(D_i) > 0 \end{cases} \quad (5.76)$$

ШАГ 3. Определяется компонента D_i с наилучшим приоритетом в смысле условия

$$R(t) = \min\{R(i) : i=1, \dots, \tau\},$$

где $\tau = k$ — количество компонент после выполненных k измерений, а также наилучшее достигнутое значение целевой функции среди допустимых измерений

$$\tilde{f}^* = \min\{\min\{f^i : g_j^i \leq 0 \forall j \in \{1, \dots, m\}\}; +\infty\}.$$

ШАГ 4. Если $R(i)$ конечно и $R(i) > \tilde{f}^* - \varepsilon$, где $\varepsilon > 0$ — заданная точность по значению целевой функции, то метод останавливается, т.к. оценка решения найдена с указанной точностью. Если $R(t) = +\infty$, то допустимая область пуста, решений нет, метод остановится. В остальных случаях — переход на шаг 5.

ШАГ 5. Компонент D_i делится по большему ребру на три равные гиперпараллелепипеда, заменяющие собой делящийся. В центрах новых компонент проводятся испытания функции Q , причем в средней из них используется уже имеющийся результат измерения. Характеристики новых компонент вычисляются по формулам (5.75), (5.76). Выполняется возврат на шаг 3.

Заметим, что способ задания характеристик в (5.76) позволяет без потери решения отбрасывать компоненты, где ограничения заведомо нарушаются. Среди остальных лучшими будут те, где нижняя оценка целевой функции в D_i меньше.

В качестве иллюстрации работы метода рассмотрим задачу (рис.5.9)

$$f(y) = y_1^2 + y_2^2 - \cos(18y_1) - \cos(18y_2) \rightarrow \min, \quad (5.81)$$

$$y_1, y_2 \in [-0.3; 0.7], \quad g_1(y) = y_1 - y_2 \leq 0.$$

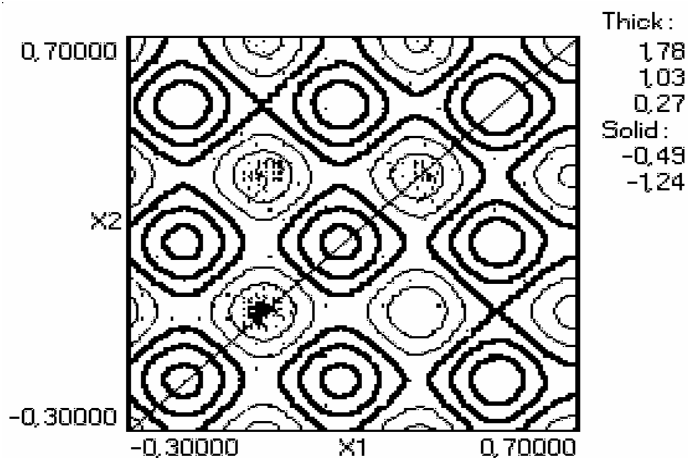


Рис.5.9. Пример размещения испытаний в методе деления на три

Укажем методу следующие значения констант Липшица $L_0=20$, $L_1=1.8$ и точность поиска $\varepsilon=0.05$. Следует обратить внимание на то, что данный метод не умеет оценивать константы класса функций и их приходится задавать. При этих параметрах метод останавливается после 379 испытаний. Картина размещения измерений, проведенных методом, показана на рисунке. Допустимой области задачи соответствует верхний из двух треугольников, на которые область D разделяется функциональным ограничением. Результат расчета показывает, что метод строит адаптивное неравномерное покрытие, чуть менее плотное в недопустимой области. При завышении значений констант L_0 , L_1 число измерений возрастет.

 **Контрольные вопросы и упражнения.**

Докажите, что при выполнении условий применимости метод деления на три остановится через конечное число шагов и либо найдет допустимую точку со значением целевой функции, отличающимся от оптимального f^ не более чем на ε , либо установит, что допустимая область Y пуста.*

5.7.2. Диагональные компонентные методы Я.Пинтера

Этот класс методов предполагает решение задач без ограничений, в которых допустимой областью является гиперпараллелепипед D из (5.3), последовательно разбиваемый в процессе решения задачи на компоненты, которые также являются гиперпараллелепипедами. Будем обозначать их как D_i^k ($i=1, \dots, \tau(k)$). Здесь k — номер итерации метода, а i — порядковый номер компоненты в общем списке, а $\tau(k)$ — текущее количество компонент. Для оценки характеристики $R(D_i^k)$ компоненты-параллелепипеда D_i^k используется его главная диагональ $[a(i); b(i)]$. В вершинах, которые эта диагональ соединяет, должны быть вычислены и запомнены значения функции. По ним определяется характеристика $R(\cdot)$, точно так же, как она определялась бы у некоторого одномерного метода на отрезке, соответствующем этой диагонали. Это можно сделать, поскольку диагональ — одномерный линейный объект. Например, для класса функций $\Phi = Lip_L(D)$ можно воспользоваться соотношениями для вычисления характеристики метода С.А. Пиявского. Затем на диагонали $[a(t(k)); b(t(k))]$ наилучшей по значениям $R(D_i^k)$ ($i=1, \dots, \tau(k)$) компоненты $D_{i(t(k))}^k$ вычисляется точка деления $\tilde{y}(t(k))$ этой диагонали по правилу вычисления точки испытания одномерного метода. Далее Я. Пинтером предлагается два способа разделения области $D_{i(t(k))}^k$ с помощью этой точки. Одна схема деления — на 2^N , а другая — пополам.

По схеме деления на 2^N через точку $\tilde{y}(t(k))$ проводятся плоскости, параллельные граням делящегося N -мерного параллелепипеда. Он распадается на 2^N новых. В концах их главных диагоналей вычисляются и запоминаются недостающие измерения функции. Точка $\tilde{y}(t(k))$ также является одной из точек измерений. В общей сложности операция разделения на 2^N требует выполнения $2 \cdot 2^N - 3$ новых измерений [40]. Это иллюстрируется левым рисунком 5.10.

При схеме деления пополам N -мерный параллелепипед $D_{i(t(k))}^k$ разделяется на два новых плоскостью, проведенной через точку $\tilde{y}(t(k))$ перпендикулярно большому ребру делимой компоненты $D_{i(t(k))}^k$. В общей сложности при этом

проводится два новых измерения, причем в $\tilde{y}(t(k))$ измерение не проводится. Это иллюстрируется правым рисунком 5.10.

После деления $D_{t(k)}^k$ процесс повторяется с обновленным набором компонент D_i^{k+1} ($i = 1, \dots, \tau(k+1)$). Остановка процесса происходит, когда объем наилучшей компоненты станет меньше заданной точности.

В теоретическом плане в смысле условий сходимости ничего принципиально нового в многомерном случае при использовании этой схемы не появляется, однако интуитивно понятные схемы дробления, предложенные Я. Пинтером, имеют существенные недостатки, исследованные в работе Я.Д. Сергеева [55].

Основная проблема заключается в следующем. Каждая из новых компонент–параллелепипедов, возникающих в процессе деления, требует наличия результатов измерений функции всего в двух точках на концах ее главной диагонали. Однако, проведение таких измерений во всех имеющихся компонентах приводит, как это было установлено в [55], к весьма неприятному эффекту, а именно, в каждой из компонент D_i^k в общей сложности окажется проведено более двух испытаний. Можно доказать (смотрите работу [55]), что при схеме деления на 2^N частей испытания функции проводятся во всех 2^N вершинах гиперпараллелепипедов с диагоналями $[a(t(k)); \tilde{y}(t(k))]$ и $[\tilde{y}(t(k)); b(t(k))]$ и, как минимум, в $(2^{N-1} + 1)$ вершинах других компонент. Для $N=2$ это иллюстрирует левый рис. 5.10. Использованные на рисунке обозначения предполагают, что одна из новых компонент (а именно компонента с наименьшими координатами левого конца диагонали) замещает в памяти делящуюся компоненту, получая ее номер $t(k)$, а остальные записываются после последней имевшейся, получая номера, следующие за ее номером $\tau(k)$.

Аналогично, при схеме деления $D_{t(k)}^k$ пополам в каждой из двух образовавшихся частей функция $f(y)$ вычисляется в общей сложности (см. [55]) в трех точках. Это иллюстрирует правый рис. 5.10.

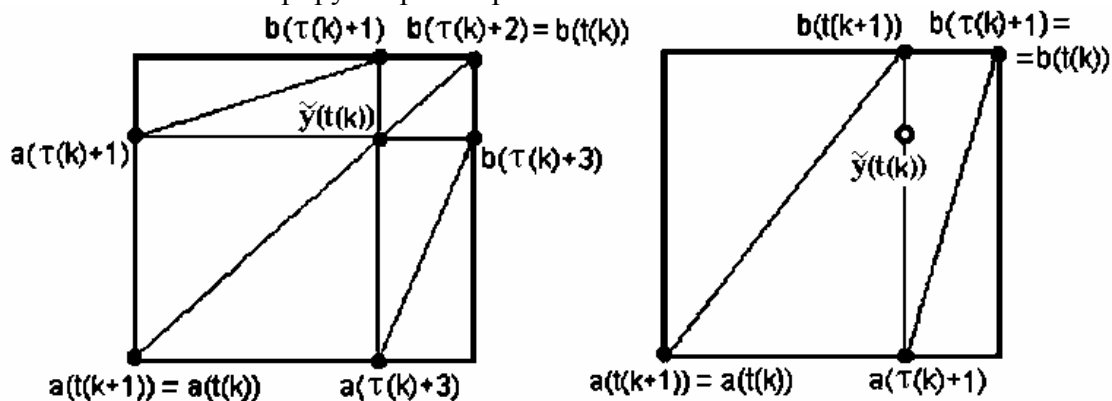


Рис. 5.10. Две схемы деления компонент в диагональном методе Я. Пинтера

Другим недостатком схем деления Я. Пинтера является (см. работу [40]) потеря информации о возможной близости некоторых вершин подобластей, возникших на разных итерациях. Примеры таких ситуаций обведены контурами на левом и правом рис.5.11. Цифрами показаны номера итераций, на которых в соответствующих точках проводились испытания $f(y)$. На правом рис. 5.11 отмечена ситуация дублирования испытания в одной из точек на 7–м шаге.

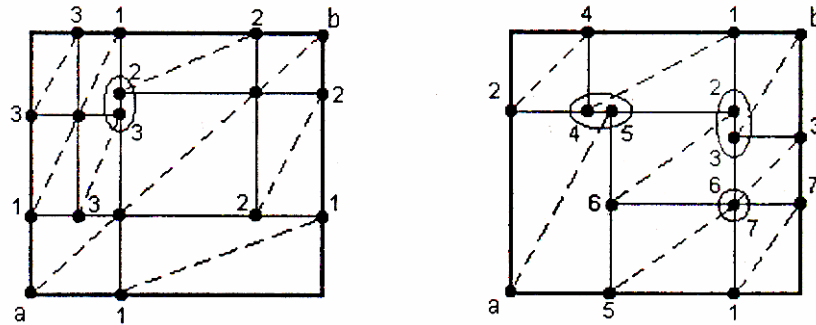


Рис. 5.11. Избыточность испытаний и потеря информации о близости точек в двух схемах дробления Я. Пинтера

5.7.3. Эффективные диагональные компонентные методы на основе адаптивных диагональных кривых

Для преодоления отмеченных выше недостатков Я.Д.Сергеевым был предложен новый для схемы Я.Пинтера способ разделения, в котором гиперпараллелепипед делится по самому длинному ребру на три части и в каждой из новых компонент строятся специальным образом согласованные диагонали для дальнейшего вычисления характеристик [40]. Способ разбиения и выбор диагоналей поясняется на рис. 5.12. При этом в [40] использован такой выбор согласования, что из диагоналей компонент-параллелепипедов D_i^{k+1} ($i = 1, \dots, \tau(k+1)$) образуется непрерывная ломаная, порождающая своеобразную адаптивную развертку области. Заметим, что непрерывность кривой обеспечивается заменой на каждом текущем шаге ее участка, являющегося диагональю делимой компоненты, на непрерывную ломаную, составленную из трех согласованно выбранных диагоналей новых компонент, размещаемых так, что концы этой ломаной совпадают с концами замещаемой ею диагонали разделяемой компоненты (рис. 5.12).

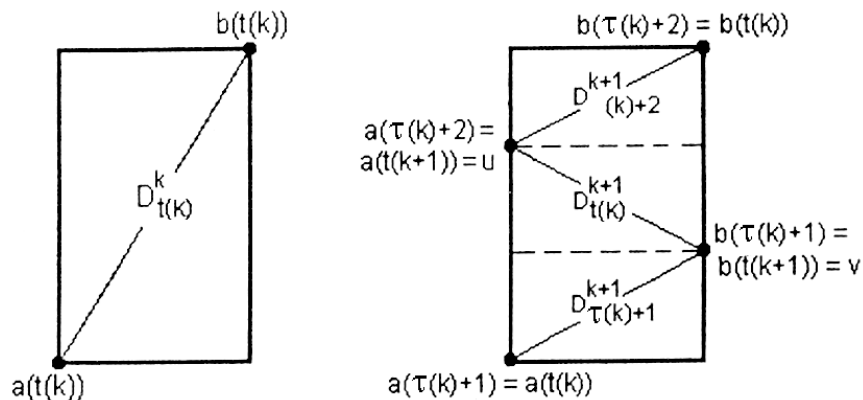


Рис. 5.12. Новая схема деления на три, замена диагонали трехзвенным участком адаптивной диагональной кривой

На рисунке 5.13 слева показано разбиение области D для $N=2$ при помощи схемы Я.Д. Сергеева, а справа — порожденная этим разбиением адаптивная диагональная кривая. На рис. 5.13 выделены точки испытаний функции и в левой части рисунка показаны номера итераций, на которых они были выполнены. Узлы возникшей адаптивной развертки маркированы в правой части рисунка символами T_i .

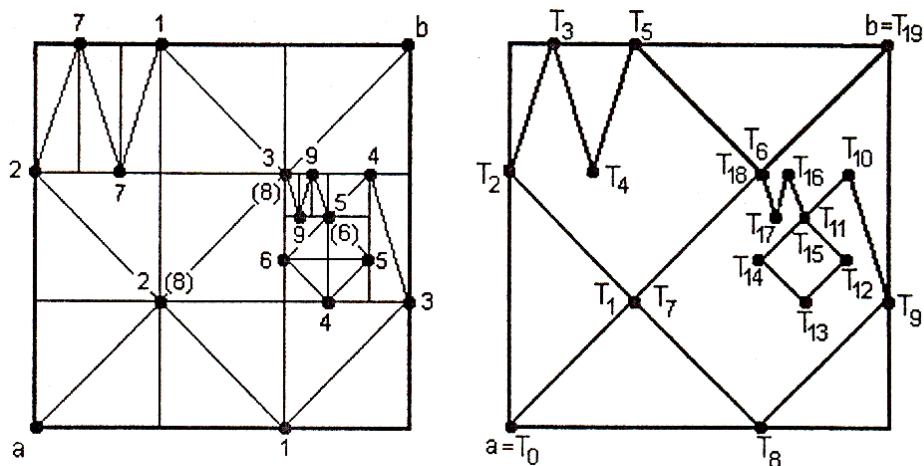


Рис. 5.13. Разбиение области по схеме Я.Д. Сергеева и вид адаптивной диагональной кривой

Поскольку кривая может иметь самопересечения (например, $T_1 = T_7$, $T_{18} = T_6$, $T_{11} = T_{16}$), то реализация метода требует организации такой системы хранения результатов выполненных итераций, при которой обеспечивается быстрый поиск и считывание из памяти результатов уже выполненных испытаний для совпадающих узлов адаптивной диагональной кривой [40].

Использование адаптивных диагональных кривых является одним из перспективных направлений в области глобальной оптимизации.

5.7.4. Компонентные методы, основанные на триангуляции области поиска

Особенностью описанных выше компонентных методов является использование компонент-параллелепипедов, а также возможность одновременного учета в каждой компоненте лишь небольшого числа испытаний, обычно — одного или двух.

Триангуляционные методы, предложенные в [19], используют адаптивное разбиение области D на многогранники S_i с $(N+1)$ -вершинами, образующими симплексы и порождающими триангуляцию области. Особенностью подхода является то, что множество вершин всех многогранников S_i совпадает с Y_k — множеством точек всех проведенных испытаний. Характеристика $R(i)$ для каждой компоненты S_i вычисляется при использовании совокупности размещенных в вершинах S_i $(N+1)$ -го испытания функции f , что позволяет более полно учитывать поисковую информацию.

Поскольку точки используемых в S_i измерений образуют N -мерные симплексы, то такие методы можно назвать симплекс-методами многоэкстремальной оптимизации или SM-методами.

SM-методы построены для двух классов целевых функций f : класса липшицевых функций $\Phi = Lip(D)$ и класса дифференцируемых функций с липшицевой производной по направлениям

$$\forall y' \neq y'' \in D \text{ и } v = (y' - y'') / \|y' - y''\|: \quad \left| \frac{\partial Q(y')}{\partial v} - \frac{\partial Q(y'')}{\partial v} \right| \leq L^* \|y' - y''\|. \quad (5.77)$$

В последнем случае характеристика компонента может учитывать результаты измерений не только $f(y)$, но и ее градиента.

ПРИНЦИПИАЛЬНАЯ СХЕМА АЛГОРИТМА

ШАГ 1. Задаются параметры метода k . Начальные испытания проводятся во всех вершинах D , его геометрическом центре и, в зависимости от размерности задачи N и выбранного вида начальной триангуляции, — в центрах граней D определенных размерностей. Строится начальное разбиение, определяющее триангуляцию D . Вершины компонент S_i симплекс-разбиения размещаются в точках проведенных испытаний. Запоминается k — количество выполненных испытаний и n_k — количество симплексов разбиения.

ШАГ 2. Для каждой из компонент S_i оценивается константа класса L , получается ее локальная оценка $l(S_i)$. Вычисляется глобальная (общая) оценка константы класса

$$l_k^* = \max\{l(S_i): i=1, \dots, n_k\}$$

и ее завышенная оценка L_k^* , лежащая в пределах $\gamma_2 * l_k^* / \gamma_1 \leq L_k^* \leq \gamma_2 * l_k^*$ при $l_k^* > 0$ и равная l при $l_k^* = 0$ ($1 < \gamma_1 < \gamma_2$). По этим характеристикам определяется локализованная оценка $L(S_i)$ константы класса, используемая для компоненты S_i при вычислении ее приоритета.

$$L(S_i) = \max\{L_*, \underline{L}_i\},$$

где $L_* > 0$ — вычисляемая по всем $l(S_i)$ заниженная общая оценка константы класса, а \underline{L}_i — скорректированное взвешенное среднее между завышенной локальной оценкой $\gamma_2 * l(S_i) / \gamma_1$ и завышенной глобальной L_k^* . Весовые коэффициенты зависят от диаметра $d(S_i)$ компоненты S_i и заданного порогового значения d^* для этого диаметра: $\underline{L}_i = L_k^*$, при $d(S_i) > d^*$, а в противном случае,

$$\underline{L}_i = (\gamma_2 * l(S_i) / \gamma_1) (1 - d(S_i) / d^*) + L_k^* (d(S_i) / d^*).$$

ШАГ 3. Для всех компонент S_i симплекс-разбиения с использованием $L(S_i)$, — локализованных оценок констант класса, вычисляются характеристики $R(i) = R(S_i)$ и определяется наиболее приоритетная компонента S_t , где $R(S_i) = \min\{R(S_i): i = 1, \dots, n_k\}$. Если диаметр $d(S_t) \leq \varepsilon_x$, то метод останавливается, иначе — переходит на шаг 4.

ШАГ 4. Точка очередного испытания y_{k+1} размещается на наибольшем ребре r_t компоненты S_t так, что делит его в пропорции ν , где $0 < \nu^* \leq \nu \leq (1 - \nu^*)$, ν^* — заданный параметр метода ($0 < \nu^* \leq 0,5$).

ШАГ 5. Проводится испытание в выбранной точке, пополняется множество точек испытаний и их результатов $Y_{k+1} = Y_k \cup \{y_{k+1}\}$, $\omega_{k+1} = \omega_k \cup \{(f, y_{k+1})\}$, компонента S_t^k разделяется на две новые с помощью точки y_{k+1} . Определяются все остальные компоненты S_i разбиения, содержащие ребро r_t (это можно сделать без полного перебора компонент). Выполняется их разделение на две части с использованием в качестве их новой вершины точки выполненного последнего испытания. Предпринимаются дополнительные действия, описанные в [19], по обеспечению условия $R_{\min}(S_i) / R_{\max}(S_i) > \beta$, где $0 < \beta < 1$ — параметр метода, а R_{\min} , R_{\max} — длины наибольшего и наименьшего ребра многогранника S_i (их отношение характеризует степень вырождения S_i). Возникает новое разбиение области D . Далее полагается $k = k + 1$, корректируется число компонент и выполняется переход на шаг 2.

Замечание.

1. В алгоритм обычно встраивается дополнительное адаптивное преобразование функции решаемой задачи, улучшающее ее свойства [19].
2. Эффективная программная реализация алгоритма требует использования достаточно сложно организованных динамических структур данных.

Приведенная принципиальная схема требует конкретизации способа вычисления локальных оценок параметров класса функций, а также способа подсчета характеристик компонент $R(i)$. Эти две части метода существенно зависят от свойств класса функций. Здесь они будут кратко описаны для липшицевых функций (класс $Lip(D)$). Читателей, интересующихся учетом измерений градиента в SM-методах при оптимизации функций из класса (5.77) адресуем к работе [20].

Итак, рассмотрим класс функций $Lip(D)$. Одну из компонент разбиения обозначим через S , ее вершины — через v^0, \dots, v^N , а результаты испытаний — через f^0, \dots, f^N . Локальную оценку константы L для S обозначим через $l(S)$, а локализованную — через $L(S)$. Способы вычисления локализованных оценок по локальным приведены в описанной выше принципиальной структуре алгоритма.

Локальная оценка определяется согласно соотношению

$$l(S) = \max \{ |f^i - f^j| / \|v^i - v^j\| : i=0, \dots, N-1; j=i+1, \dots, N \}.$$

Значение функции приоритета будем строить как оценку минимума миноранты

$$f^-_s(y) = \max \{ f^j - L(S) \|v^j - y\| : v^j \in S \}$$

функции f в S , по результатам испытаний f в вершинах S . Для вычисления характеристики $R(S)$ компоненты S вводится функция

$$\Psi(v, C, y) = C + L(S) \|v - y\|,$$

где C и y определяются из системы нелинейных уравнений

$$\Psi(v^j, C, y) = f^j \quad (j=0, \dots, N). \quad (5.78)$$

В [19] описан вычислительный метод решения (5.78), основанный на сведении к линейной системе. Решения обозначим через $C(S)$ и $y(S)$. В [19] доказано, что для класса $Lip(D)$ при $y(S) \in S$ выполняется: $C(S) = \min \{ f^-_s(y) : y \in S \}$. Эта ситуация показана на рис.5.14 слева.

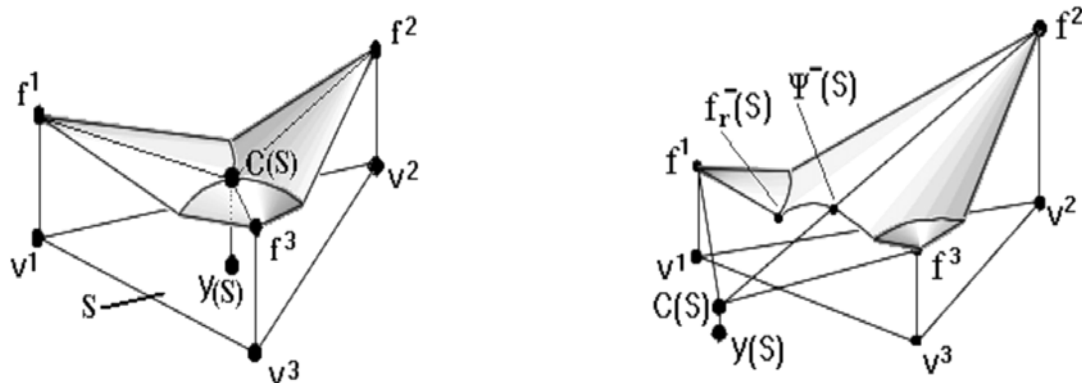


Рис.5.14. Нижние оценки функции по измерениям в вершинах симплекса S

Если же $y(S) \notin S$, то для уже вычисленных $C(S)$ и $y(S)$ методами квадратичного программирования определяется значение $\Psi^-(S) = \min \{ \Psi(v, C(S), y(S)) : v \in S \}$, а

также определяется $f_r^-(S)$ — минимум миноранты $f_s^-(y)$ на ребрах S .
 Окончательно получим следующее правило вычисления характеристики
 компоненты

$$R(S) = \begin{cases} C(S), & y(S) \in S \\ \min\{\Psi^-(S), f_r^-(S)\}, & y(S) \notin S. \end{cases} \quad (5.79)$$

Геометрическая иллюстрация для двух случаев приведена на рис.5.14.

Вычислительные иллюстрации

Приведем несколько вычислительных иллюстраций. На рис.5.15 приведены
 картины размещения точек испытаний при решении следующих двух задач.
 Левому рисунку соответствует задача

$$f(y) = -(\sin(4y_1+1)+2\sin(6y_2+2)) \rightarrow \min, \quad (5.80)$$

$$y_1, y_2 \in [-2; 2].$$

решенная SM–методом при следующих параметрах $\gamma_1=1.2, \gamma_2=1.8, d^*=0.1 \text{ diam}(D),$
 $\varepsilon_x=0.01$.

На правом рисунке показано размещение точек при решении задачи

$$f(y) = y_1^2 + y_2^2 - \cos(18y_1) - \cos(18y_2) \rightarrow \min, \quad (5.81)$$

$$y_1, y_2 \in [-0.3; 0.7].$$

при $\gamma_1=1.2, \gamma_2=1.8, d^*=0.1$ от диаметра области $D, \varepsilon_x=0.005$.

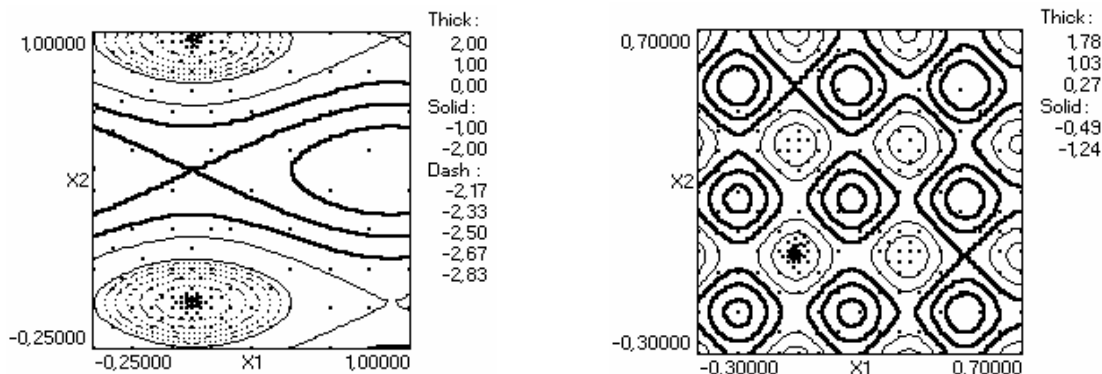


Рис.5.15. Размещение точек испытаний SM–методом

На рис.5.16 показана триангуляция области D , построенная SM–методом на
 некотором промежуточном шаге в процессе решения задачи (5.80) с двумя
 глобальными минимумами.

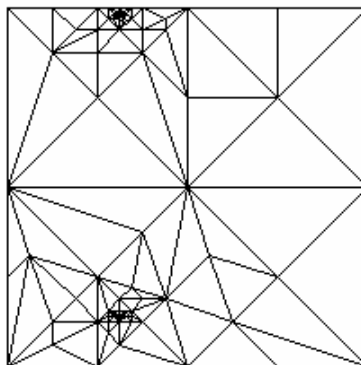


Рис.5.16. Триангуляция области, построенная SM–методом в задаче (5.80)

Эта иллюстрация работы метода соответствует размещению измерений, представленному слева на рис.5.15, но при меньшем числе проведенных измерений.

На рисунке видно, как процесс триангуляции адаптируется к структуре решаемой задачи.

Лист регистрации изменений

Дата	Автор	Комментарии
??.06.03	Гришагин В.А.	Первоначальная версия главы 5
??.07.03	Гришагин В.А.	Правка текста
26.05.03	Городецкий С.Ю.	Создание копии раздела 5.7 с использованием версии раздела 5.7, написанной В.А.Гришагиным
27.08.03	Городецкий С.Ю.	Начата переработка раздела 5.7
27.08.03– 15.09.03	Городецкий С.Ю.	Основная переработка раздела 5.7
20.09.03	Городецкий С.Ю.	Подготовка рисунков в раздел 5.7
12.10.03– 15.10.03	Городецкий С.Ю.	Корректурa раздела 5.7 и его расширение
18.10.03	Городецкий С.Ю.	Окончательная корректурa 5.7
09.10.03– 17.10.03	Гришагин В.А.	Изменение обозначений в разделах 5.1–5.6, доработка, корректурa
18.10.03	Городецкий С.Ю.	Присоединение новых версий раздела 5.7 к главе 5
19.10.03	Городецкий С.Ю.	Изменение стилей оформления в разделах 5.1–5.6

Глава 6. Методы построения оценок множества слабо эффективных точек, не использующие параметрических сверток

В разделе 1.3 второй главы были введены понятия решений по Парето и Слейтеру (см. определение 1.4) для задач (1.21)–(1.23) с векторным критерием.

Напомним, что в пунктах 1.3.2–1.3.6 рассматривались различные схемы компромисса и было показано, что их использование приводит к параметрической скаляризации векторной задачи с использованием некоторой функции свертки. При заданных параметрах используемого метода скаляризации из решения возникающей вспомогательной оптимизационной задачи находится одно из эффективных y^* или слабо эффективных y^o решений исходной задачи. Изменяя значения параметров, можно последовательно (поочередно) определять оценки различных решений из множеств Y^* или Y^o .

Все методы параметрической скаляризации имеют тот общий недостаток, что при их использовании остается открытым вопрос об оценивании всего множества эффективных (слабо эффективных) решений в целом. В пункте 1.3.7 была поставлена задача такого оценивания. В данной лекции мы вернемся к этому вопросу и рассмотрим несколько подходов и методов, альтернативных к ранее изученным процедурам параметрической скаляризации.

Материал лекции, в основном, опирается на результаты, приведенные в работах Евтушенко Ю.Г. и Потапова М.А. (см. [24], а также более подробный вариант изложения в [25]), работах Городецкого С.Ю. [17,18], а также Маркина Д.Л. и Стронгина Р.Г. [30].

Далее будет описано несколько близких подходов, представленных в этих работах:

- метод неравномерных покрытий Евтушенко Ю.Г. и Потапова М.А., позволяющий строить ε –оптимальные оценки всего множества эффективных точек Y^* ;
- метод построения перестраиваемого семейства непараметрических сверток, приводящих к построению ε –оптимальной оценки множества слабо эффективных точек Y^o , а также использующее эти свертки семейство одношагово–оптимальных методов для оценивания множества Y^o в целом (результат Городецкого С.Ю.);
- метод построения непараметрической свертки, порождающей скалярную задачу, множество точек глобальных минимумов которой совпадает с множеством слабо эффективных точек исходной задачи (результат Маркина Д.Л., Стронгина Р.Г.).

Все результаты будут представлены с общей позиции, часто не совпадающей с авторской версией изложения.

Поскольку лекция направлена, в первую очередь, на изучение перспективных методов непараметрической скаляризации, исходная задача (1.21)–(1.23) будет рассматриваться в упрощенной постановке, без функциональных ограничений

$$f(y) = (f_1(y), \dots, f_n(y)) \rightarrow \min, \quad y \in Y = D = \{y \in R^N : a \leq y \leq b\}. \quad (6.1)$$

Заметим, что учет функциональных ограничений не вызывает принципиальных трудностей, но делает изложение более громоздким.

Напомним некоторые из ранее использованных обозначений, а также добавим новые. Пусть $F = f(Y) = \{z = f(y) : y \in Y\}$ — множество возможных векторных оценок, $Y_k = \{y^j : y^j \in Y (j = 1, \dots, k)\}$ — множество точек выполненных измерений вектор-функции f , а $F_k = f(Y_k) = \{z^j = f(y^j) : y^j \in Y_k\}$ — множество результатов этих измерений. В пункте 1.3.2. были использованы операторы P и S , определенные на множествах оценок и выделяющие из них подмножества Парето и Слейтера. При этом $P(F)$ и $S(F)$ — будут множествами Парето и Слейтера исходной задачи, а их прообразы при отображении f (обозначаемые через Y^* и Y^o) — множествами эффективных и слабо эффективных точек в пространстве параметров. Кроме того, $P(F_k)$, $S(F_k)$ — будут аппроксимациями множеств Парето и Слейтера, построенными по конечному набору вычисленных векторных оценок F_k , а прообразы этих множеств, порождаемые какой-либо однозначной ветвью отображения f^{-1} (обозначим их через Y_k^* , Y_k^o), — будут текущими оценками множеств Y^* , Y^o .

Приведем алгоритм построения таких оценок на примере Y_k^o .

РЕКУРРЕНТНЫЙ АЛГОРИТМ ПОСТРОЕНИЯ ОЦЕНОЧНЫХ МНОЖЕСТВ Y_k^o .

ШАГ 0. Вначале полагается $Y_k^o = \{y^1\}$, $k = 1$

ШАГ 1. Появляется результат очередного измерения f^{k+1} , выполненного в точке y^{k+1} .

ШАГ 2. Если $\exists y^j \in Y_k^o$ ($0 < j \leq k$), что $f^{k+1} < f^j$, то все такие точки y^j исключаются из Y_k^o , а точка y^{k+1} включается в Y_k^o , затем полагается $Y_{k+1}^o = Y_k^o$. Иначе, если $\exists y^j \in Y_k^o$ ($0 < j \leq k$), что $f^j < f^{k+1}$, то полагается $Y_{k+1}^o = Y_k^o$, иначе y^{k+1} добавляется к Y_k^o , т.е. $Y_{k+1}^o = Y_k^o \cup \{y^{k+1}\}$.

ШАГ 3. Принимается $k := k+1$ и выполняется переход на шаг 1.

6.1. Основные принципы непараметрической скаляризации

В пункте 1.3.7 было дано определение 1.6 множества Y_ε^* — ε -оптимальных (по Парето) решений задачи. Приведем его незначительную модификацию, необходимую для дальнейшего изложения, и дополним определением ε -оптимального решения по Слейтеру.

Определение 6.1. Пусть $e \in R^n$, $e > 0$ и $\|e\| = 1$, $\varepsilon > 0$. Множество $Y_\varepsilon^* \subseteq Y$ будем называть ε -оптимальным по Парето решением задачи (6.1), если $\forall y_\varepsilon^* \in Y_\varepsilon^* \exists y_\varepsilon \in Y_\varepsilon^*$, что $f(y_\varepsilon^*) \leq f(y_\varepsilon) + \varepsilon e$ и в Y_ε^* нет двух разных точек y_ε^{*1} , y_ε^{*2} , что $f(y_\varepsilon^{*1}) \leq f(y_\varepsilon^{*2})$.

Множество $Y_\varepsilon^o \subseteq Y$ будем называть ε -оптимальным по Слейтеру решением задачи (6.1), если $\forall y_\varepsilon^o \in Y_\varepsilon^o \exists y_\varepsilon^o \in Y_\varepsilon^o$, что $f(y_\varepsilon^o) \leq f(y_\varepsilon^o) + \varepsilon e$, и в Y_ε^o нет двух разных точек $y_\varepsilon^{o1} \neq y_\varepsilon^{o2}$, что $f(y_\varepsilon^{o1}) = f(y_\varepsilon^{o2})$ или $f(y_\varepsilon^{o1}) < f(y_\varepsilon^{o2})$.

6.1.1. Метод сведения к скалярной задаче с перестраиваемой целевой функцией

Рассмотрим идею построения ε -оптимального решения, следуя работам [24, 25].

Пусть $\tilde{Z} \subset R^n$ — некоторое замкнутое множество точек в пространстве оценок. Введем специальную функцию

$$\rho(z, \tilde{Z}) = \min_{\tilde{z} \in \tilde{Z}} \max_{i=1, \dots, n} ((\tilde{z}_i - z_i) / e_i). \quad (6.2)$$

Ей можно дать следующую геометрическую интерпретацию. Рассмотрим введенное в разделе 1.3.5. множество

$$R^+(z') = \{ z \in R^n : z_i \geq z'_i, (i=1, \dots, n) \}$$

и построим еще одно вспомогательное множество в виде объединения

$$(\tilde{Z})^+ = \bigcup_{z \in \tilde{Z}} R^+(z).$$

Его границу обозначим $\partial(\tilde{Z})^+$. Очевидно, геометрическое место точек со значением $\rho(z, \tilde{Z}) = 0$ будет совпадать с этой границей. Будем теперь смещать множество $\partial(\tilde{Z})^+$ в направлении $(-e)$ на величину C (смещение с $C < 0$ трактуется как смещение в направлении $+e$). Поверхность $\partial(\tilde{Z})^+ - Ce$, полученная из $\partial(\tilde{Z})^+$ параллельным сдвигом на вектор $-Ce$, будет совпадать с геометрическим местом точек, где $\rho(z, \tilde{Z}) = C$ (рис. 6.1).

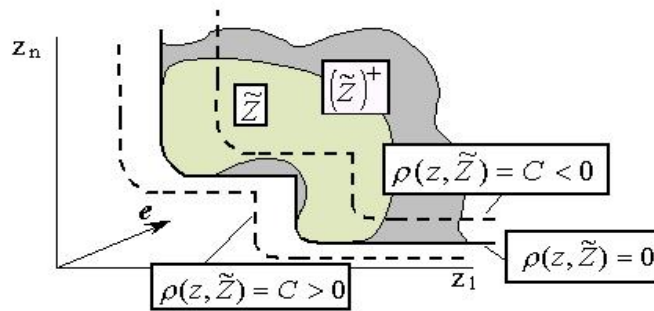


Рис.6.1. Структура изолиний функции $\rho(z, \tilde{Z})$

Используем теперь в качестве \tilde{Z} конечное слейтеровское подмножество $S(F_k)$ результатов выполненных измерений. Рассмотрим семейство поверхностей

$$\partial_C = \{ z \in R^n : \rho(z, S(F_k)) = C \}, \quad (6.3)$$

показанных на рис.6.2.

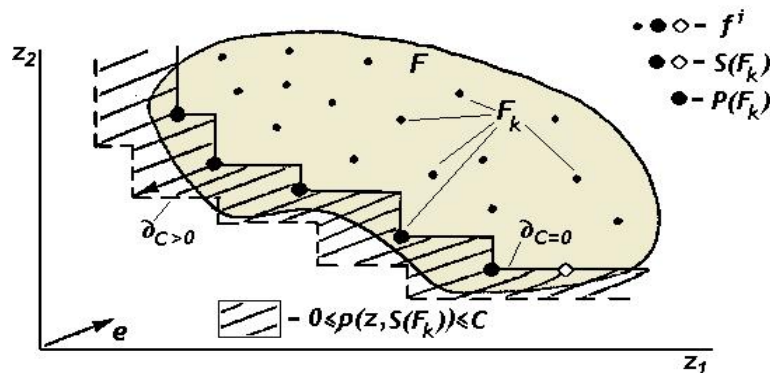



Рис.6.2. Вид оценок $P(F_k)$ и $S(F_k)$ и область «улучшения» на величину от 0 до C

Имеют место несколько очевидных свойств.

Свойство 6.1. Пусть $z^j \in F_k$. При этом $z^j \in S(F_k)$ тогда и только тогда, когда $\rho(z, S(F_k))=0$, т.е. $z^j \in \partial_{C=0}$.

Свойство 6.2. Если для точки $y \in Y$ значение $\rho(f(y), S(F_k))=C > 0$, то включение в F_k дополнительного значения $z=f(y)$ приведет к смещению участка поверхности $\partial_{C=0}$ на величину C вдоль направления $(-e)$.

Описанная в свойстве 6.2 ситуация эквивалентна соответствующему улучшению текущей оценки множества Слейтера $S(F_k)$.

 **Замечание.** Функцию $\rho(z, S(F_k))$ можно трактовать как меру улучшения оценки множества Слейтера $S(F_k)$ за счет учета результата $z=f(y)$ очередного измерения векторного критерия.

Введем непараметрическую функцию свертки

$$\Psi_{F_k}(z) = -\rho(z, S(F_k)), \quad (6.4)$$

уменьшение значений которой эквивалентно соответствующему увеличению значения меры $\rho(\cdot)$.

Приведенное замечание, а также свойство 6.2 могут служить основанием для перехода к следующему перестраиваемому семейству задач.

$$\begin{aligned} \Psi_{F_k}(f(y)) \rightarrow \min, y \in Y = D, \\ F_{k+1} = F_k \cup \{z^{k+1} = f(y^{k+1})\}, k = 1, 2, \dots, \end{aligned} \quad (6.5)$$

где y^{k+1} — точка очередного измерения векторного критерия f при решении перестраиваемой задачи (6.5).

Лемма 6.1. [17, 18] Для того, чтобы множество Y_k^o , полученное по рекуррентному алгоритму построения оценочных множеств (см. алгоритм перед разделом 6.1), являлось ε -оптимальным (по Слейтеру) решением задачи (6.1) необходимо и достаточно, чтобы $\forall z^o \in S(F)$

$$0 \leq \rho(z^o, f(Y_k^o)) \leq \varepsilon, \quad (6.6)$$

где $f(Y_k^o) \equiv S(F_k)$.

НЕОБХОДИМОСТЬ. Пусть Y_k^o — ε -оптимальное по Слейтеру решение. Возьмем произвольное $z^o \in S(F)$. При этом найдутся $y^o \in Y^o$, что $z^o = f(y^o)$. По определению 6.1, для $y^o \in Y^o \exists y_k^o \in Y_k^o$, что $f(y_k^o) \leq f(y^o) + \varepsilon e$ и, следовательно,

$$\forall i = 1, \dots, n \quad (f_i(y_k^o) - f_i(y^o)) / e_i \leq \varepsilon.$$

Отсюда следует, что

$$\rho(z^o, f(Y_k^o)) = \min_{y \in Y_k^o} \max_{1 \leq i \leq n} ((f_i(y) - f_i(y^o)) / e_i) \leq \varepsilon.$$

Далее, поскольку для $z^o = f(y^o) \in S(F)$ не существует точки $y \in Y$, что $f(y) < f(y^o)$, то $\forall y \in Y$ выполнится неравенство

$$\max_{1 \leq i \leq n} ((f_i(y) - f_i(y^o)) / e_i) \geq 0.$$

Такой же знак будет у минимального значения этой функции, найденного по всем $y \in Y_k^o$. Используя вид функции $\rho(\cdot)$ из (6.2), отсюда видим, что для $z^o = f(y^o)$ $\rho(z, f(Y_k)) \geq 0$. Необходимость доказана.

ДОСТАТОЧНОСТЬ. Пусть $\forall z^o \in S(F)$ выполняется неравенство (6.6). Тогда $\forall y^o \in Y^o$

$$\min_{y_k^o \in Y_k^o} \max_{1 \leq i \leq n} ((f_i(y_k^o) - f_i(y^o)) / e_i) \leq \varepsilon.$$

Следовательно, $\exists y_k^o \in Y_k^o$, что $f(y_k^o) \leq f(y^o) + \varepsilon e$. Таким образом, выполнено главное свойство ε -оптимального (по Слейтеру) решения. Остальные свойства обеспечиваются алгоритмом построения множества Y_k^o . Таким образом, Y_k^o — ε -оптимально по Слейтеру. Достаточность доказана.

Следствие 6.1.1. Если при решении задач (6.5) на k -м шаге выполнится условие $\min \{ \Psi_{F_k}(f(y)) : y \in D \} \geq -\varepsilon$, то множество Y_k^o , соответствующее оценке $S(F_k)$, является ε -оптимальным по Слейтеру решением многокритериальной задачи, представленной в (6.1).

Одношагово-оптимальные методы решения семейства задач (6.5), предложенные в работах [17, 18], будут рассмотрены позднее в разделе 6.3.

6.1.2. Метод неравномерных покрытий Ю.Г. Евтушенко и М.А. Потапова

Сейчас обратим внимание, что возможно иное использование меры $\rho(\cdot)$, применительно к классу задач с липшицевыми компонентами вектор-функций f .

Предположим, что $\forall i=1, \dots, n$ функции $f_i(y)$ удовлетворяют на D условию Липшица в метрике $\|\cdot\|$ пространства R^n с константами $L_i > 0$ (т.е. $f \in \Phi = Lip_L(D)$). Пусть

$$L = \| (L_1, \dots, L_n) \|,$$

где использована норма пространства оценок R^n . Выберем в (6.2) вектор e так, чтобы $\forall i=1, \dots, n$ $e_i = L_i / L$. Тогда нормированные функции $f_i(y) / e_i$ ($i=1, \dots, n$) будут иметь одинаковое значение константы Липшица, равные L .

Свойство 6.3. Пусть функция $f \in \Phi = Lip_L(D)$, выполнено k ее измерений и задано $\varepsilon > 0$. Тогда, если $\Delta_k^j = -\rho(f(y^j), S(F_k))$, то $\Delta_k^j \geq 0$ и $\forall y \in D$, такого что

$$\|y - y^j\| \leq r_k^j = (\Delta_k^j + \varepsilon) / L, \quad (6.7)$$

выполняется $\rho(f(y), S(F_k)) \leq \varepsilon$.

ДОКАЗАТЕЛЬСТВО непосредственно вытекает из формулы (6.2), свойств липшицевых функций и леммы 6.1.

Свойство 6.3 позволяет сокращать область поиска, исключая из D подобласти (6.7), измерения в которых не могут улучшить текущую оценку $S(F_k)$ более чем на ε по мере $\rho(\cdot)$.

В работе Евтушенко Ю.Г. и Потапова М.А. [24] доказана следующая теорема, вытекающая из свойства 6.3.

Теорема 6.1 (о неравномерном покрытии). Пусть $f \in \Phi = \text{Lip}_L(D)$. Если множество Y_k таково, что $D \subset \bigcup_{j=1}^k O_{r_k^j}(y^j)$, где радиусы шаров определяются формулой (6.7), то множество Y_k^* будет являться ε -оптимальным (по Парето) решением задачи (6.1).

Таким образом, поиск ε -оптимального решения, оценивающего в целом множество Парето, сводится к построению специального неравномерного покрытия Y_k допустимого множества D . Предложенный авторами данного подхода метод построения покрытий, удовлетворяющих теореме 6.1, основан на модификации метода деления на три (см. пункт 5.7.1 пятой главы), разработанного для однокритериальных задач и аналогичного методу, рассмотренному в работе Евтушенко Ю.Г. и Ратькина В.А. [26]. Идея метода, применительно к многокритериальной задаче, будет кратко описана в разделе 6.2.


6.1.3. Точное сведение многокритериальной задачи к скалярной с помощью свертки Д. Л. Маркина, Р. Г. Стронгина

Заметим, что в целях сохранения единого подхода к изложению материала, здесь будет использована иная мотивация рассматриваемого метода, нежели приведенная в работе авторов [30]. Итак, если в свертке (6.4), (6.2) использовать вместо $S(F_k)$ все множество возможных оценок $F=f(D)$, а также внести знак «минус» под операции \min/\max в (6.2), то придем к специальной свертке вида

$$\psi(z) = \max_{\tilde{y} \in D} \min_{1 \leq i \leq n} ((z_i - f_i(\tilde{y})) / e_i) \equiv -\rho(z, F) \quad (6.8)$$

Заменяя в свойстве 6.1 конечное множество $S(F_k)$ на F , нетрудно увидеть, что $\psi(f(y))=0$ для $y \in D$ тогда и только тогда, когда $y \in Y^0$. Кроме того, для допустимых значений y будет выполняться неравенство $\psi(f(y)) \geq 0$. Это показывает, что множество слабо эффективных точек Y^0 можно найти, решая скалярную задачу вида

$$\psi(f(y)) \rightarrow \min, \quad y \in Y=D \quad (6.9)$$

 **Замечание.** Многокритериальная задача (6.1) с помощью свертки (6.8) сводится к решению в пространстве размерности $2N$ минимаксной задачи вида

$$\min_{y \in D} \max_{\tilde{y} \in D} \min_{1 \leq i \leq n} ((f_i(y) - f_i(\tilde{y})) / e_i) \quad (6.10)$$

Этот результат приведен в работе Маркина Д.Л. Стронгина Р.Г. [30] в виде следующей теоремы.

Теорема 6.2. Множество слабо эффективных решений задачи (6.1) совпадает с подмножеством точек D , в которых функция $\psi(f(y))$ достигает глобально-минимального значения равно нулю, т.е.

$$Y^0 = \{y \in D: \Psi(f(y)) = \Psi^* = 0\} \quad (6.11)$$

Доказательство. В силу определения функции свертки в (6.8), полагая в правой части $\tilde{y} = y$, получим оценку $\Psi(f(y)) \geq \min_{1 \leq i \leq n} ((f_i(y) - f_i(y)) / e_i) = 0$.

Пусть $y^0 \in Y^0$. Тогда, в силу недоминируемости этой точки, $\forall y \in D \exists i$, что $f_i(y^0) \leq f_i(y)$. Следовательно, $\forall y \in D \min_{1 \leq i \leq n} ((f_i(y^0) - f_i(y)) / e_i) \leq 0$ и поэтому $\Psi(f(y^0)) \leq 0$. Однако отрицательные значения невозможны, значит для $y^0 \in Y^0$

$\Psi(f(y^o)) = 0$, откуда вытекает, что точки $y^o \in Y^o$ являются глобальными минимумами в задаче (6.9).

Наоборот, для $\forall y \in D \setminus Y^o$ найдется доминирующая точка $y' \in D$, что $f(y') < f(y)$, следовательно $\psi(f(y)) > 0 \quad \forall y \in D \setminus Y^o$. Таким образом, функция свертки (6.8) неотрицательна и обращается в нуль на D тогда и только тогда, когда $y \in Y^o$. Теорема доказана.

Тем самым показано точное сведение многокритериальной задачи (6.1) к скалярной минимаксной (6.10). Предложенный авторами [30] метод решения задачи (6.10) и его свойства приведены в разделе 6.4.

6.2. Реализация метода неравномерных покрытий Ю.Г. Евтушенко по схеме деления на три

Приведем принципиальную схему алгоритма, реализующего метод построения ε -оптимального решения задачи (6.1), основанный на неравномерных покрытиях. Алгоритм использует подход «divide-the-best» деления области на компоненты по схеме деления на три, рассмотренный в пункте 5.7.1. Область D будет разделяться на параллелепипеды D_y с гранями, параллельными координатным плоскостям, каждый из которых характеризуется своим геометрическим центром y , размерами и значением вектор-функции f в его центре. Будем описывать этот набор данных следующей совокупностью $\langle D_y, f_y = f(y) \rangle$

ОПИСАНИЕ АЛГОРИТМА

ШАГ 0. Задаются константы Липшица L_1, \dots, L_n , $\varepsilon > 0$, вычисляется $L = \|(L_1, \dots, L_n)\|$ и $e_i = L_i/L$ ($i=1, \dots, n$).

ШАГ 1. Проводится вычисление $f^l = f(y^l)$ в точке y^l , размещаемой в центре области D , полагается $Y_l = \{y^l\}$, $F_l = \{f^l\}$, $Y_l^* = Y_l$. Формируется начальный список Π подобластей-параллелепипедов, состоящий из одного исходного

$$\Pi = \{ \langle D_{y=y^l} = D, f_y = f(y) = f(y^l) \rangle \}.$$

В качестве текущей «лучшей» подобласти D_{y^*} выбирается единственная существующая подобласть $D_{y^*} = D_{y^l}$.

ШАГ 2. Выполняется исключение $\Pi := \Pi \setminus \{ \langle D_{y^*}, f(y^*) \rangle \}$. Подобласть D_{y^*} по большему ребру разбивается на три равных параллелепипеда $\bar{D}_{y'}$, \bar{D}_{y^*} , $\bar{D}_{y''}$ с центрами в точках y' , y^* , y'' . Проводятся измерения $f(y')$ и $f(y'')$.

ШАГ 3. Выполняется добавление в список

$$\Pi := \Pi \cup \{ \langle \bar{D}_{y'}, f(y') \rangle; \langle \bar{D}_{y''}, f(y'') \rangle \}$$

и коррекция множеств $F_k = F_k \cup \{f(y'); f(y'')\}$ и обновление Y_{k+2}^o , соответствующего $P(F_{k+2})$, полагается $k := k+2$.

ШАГ 4. Для каждого элемента списка $\langle D_y, f(y) \rangle \in \Pi$ вычисляется оценка

$$\Psi_{D_y} = \rho(f(y) - 0.5 e L \text{diam}(D_y), S(F_k)) \quad (6.13)$$

Если $\Psi_{D_y} \leq \varepsilon$, то элемент $\langle D_y, f(y) \rangle$ исключается из списка Π . Если $\Psi_{D_y} > \varepsilon$, то эта оценка запоминается для элемента и он на шаге 5 участвует в определении новой подобласти D_{y^*} с наибольшим значением Ψ_{D_y} .

ШАГ 5. Если список $\Pi = \emptyset$, выполняется останов и множество Y_k^* сообщается в качестве ε -оптимального решения. Если $\Pi \neq \emptyset$, то для выделенной «лучшей» подобласти D_{y^*} выполняется переход на Шаг 2.

Теорема 6.3. *Метод деления на три остановится через конечное число шагов, и в момент останова множество Y_k^* является ε -оптимальным (по Парето) решением задачи (6.1).*

ДОКАЗАТЕЛЬСТВО [24, 26], почти очевидно из построения метода.

6.3. Одношагово-оптимальный метод многокритериальной оптимизации на основе адаптивных стохастических моделей

Получим метод решения задачи (6.5) для класса $\Phi = Lip_L(D)$ целевых вектор-функций, удовлетворяющих условию Липшица с константами $L = (L_1, \dots, L_n)$, следуя [17, 18]. Заметим, что в силу выражения для свертки (6.4), решение задач (6.5) эквивалентно максимизации функции $\rho(f(y), S(F_k))$:

$$\rho(f(y), S(F_k)) \rightarrow \max, y \in Y = D, \quad (6.14)$$

$$F_{k+1} = F_k \cup \{z^{k+1} = f(y^{k+1})\}.$$

Начнем с построения простой адаптивной стохастической модели, учитывающей, что $f \in \Phi = Lip_L(D)$. Пусть выполнено k измерений функции $f(y)$ на множестве точек Y_k и имеется множество результатов измерений F_k . Таким образом, мы располагаем поисковой информацией $\omega_k = \{(y^j, f^j) : j = 1, \dots, k\}$.

В силу условий Липшица, для компонент векторного критерия существуют верхние и нижние оценки $f_k^\pm(y)$, построенные по поисковой информации ω_k . Их вид приведен в пункте 1.4.2.1. первой главы (формулы (1.53), (1.54)). Таким образом, покомпонентно выполняются неравенства

$$\forall y \in Y : f_k^-(y) \leq f(y) \leq f_k^+(y).$$

В пространстве значений критериев построим область

$$P_k(y) = \{z \in R^n : f_k^-(y) \leq z \leq f_k^+(y)\}. \quad (6.15)$$

Тогда $\forall y \in Y$ выполнится: $z = f(y) \in P_k(y)$. Это единственная информация о значении $f(y)$, которая заключена в результатах измерений ω_k .

Введем теперь дополнительные предположения вероятностного характера. Неизвестное значение $z = f(y)$ будем считать реализацией финитной векторной случайной величины ξ_{y^k} , распределенной в области $P_k(y)$ с некоторой заданной плотностью $P(z / \omega_k, y)$, сосредоточенной на $P_k(y)$. Конкретные корреляционные свойства $\xi_{y^k}^{y'}$ и $\xi_{y^k}^{y''}$ для различных y' и y'' определять не будем. Плотность $P(z / \omega_k, y)$ при отсутствии конкретных априорных представлений о поведении критериев можно принять постоянной на $P_k(y)$, соответствующей равномерному распределению. При наличии существенной корреляции между значениями компонент векторного критерия плотность $P(z / \omega_k, y)$ можно выбрать сосредоточенной вдоль главной диагонали области $P_k(y)$.

Совокупность принятых выше предположений о функции f , будем называть ее *простой адаптивной стохастической моделью* (рис.6.3). Модели такого типа для многоэкстремальной оптимизации были введены в пункте 1.4.2 первой главы.

Построим на основе принятой модели задачи решающее правило для выбора точки очередного измерения критериев. Для этого заметим, что вычисление нового значения $z^{k+1} = f(y^{k+1})$ приведет к существенному (большему $\delta > 0$) уточнению оценки множества Слейтера $S(F_k)$ только в том случае, когда значение

z^{k+1} попадет в область, где значение $\rho(z, S(F_k)) > \delta$. Заметим также, что согласно правилу (6.14) цель проведения новых испытаний состоит в максимизации функции $\rho(f(y), S(F_k))$.

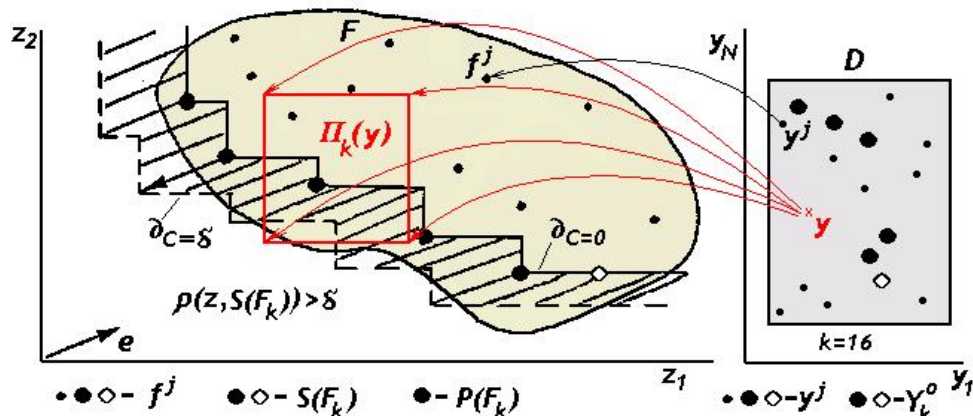


Рис.6. 3. Область неопределенности $\Pi_k(y)$ значения векторного критерия $f(y)$ в точке y и текущая оценка множества решений

Наибольшие ее значения на множестве точек проведенных испытаний равны нулю и достигаются на подмножестве точек $y=y^j \in Y_k^0$, где $f(Y_k^0)=S(F_k)$ (в пространстве критериев на рис.6.3 множеству со значением $\rho(z, S(F_k))=0$ соответствует поверхность $\partial_{C=0}$). Поэтому неравенство $\rho(z, S(F_k)) > \delta$ эквивалентно требованию улучшения текущего достигнутого (при решении задачи (6.14)) значения ее целевой функции по крайней мере на величину $\delta > 0$.

Согласно принципу одношаговой оптимальности в среднем, описанному в разделе 4.2 четвертой главы, введем функцию выигрыша на шаге (здесь использованы обозначения, несколько отличающиеся от раздела 4.2), определяющую величину выигрыша при получении результата измерения $z = f(y)$:

$$V_{k+1}(\omega_k, z) = V(\omega_k, z, \delta_k) = \begin{cases} 1, & \rho(z, S(F_k)) > \delta_k \\ 0, & \rho(z, S(F_k)) \leq \delta_k. \end{cases} \quad (6.16)$$

Тогда соответствующая ей функция среднего выигрыша на шаге

$$W_{k+1}(y, \delta_k) = \int_{\Pi_k(y)} V(\omega_k, z, \delta_k) P(z/\omega_k, y) dz = P(\rho(z = f(y), S(F_k)) > \delta_k) \quad (6.17)$$

будет определять вероятность того, что при измерении в точке $y \in Y$ в задаче (6.14) будет достигнуто улучшение текущей оценки решения по крайней мере на величину δ_k . Очередное измерение должно проводиться в такой точке y^{k+1} , для которой эта вероятность максимальна, т.е.

$$y^{k+1} = \arg \max_{y \in Y=D} W_{k+1}(y, \delta_k). \quad (6.18)$$

Заметим, что значение вероятности из (6.17) наиболее просто вычисляется в том случае, когда плотность $P(z/\omega_k, y)$ соответствует равномерному распределению, сосредоточенному на главной диагонали области $\Pi_k(y)$. В этом случае (6.17) приобретает вид

$$W_{k+1}(y, \delta_k) = \max\{\rho(f_k^-(y), S(F_k)) - \delta_k; 0\} / (\rho(f_k^-(y), S(F_k)) - \rho(f_k^+(y), S(F_k))).$$

Для завершения описания принципиальной схемы алгоритма осталось конкретизировать способ настройки параметра δ_k и определить критерий останова.

Пусть $\varepsilon \geq 0$ — заданная точность решения задачи.

МЕТОД 1. (С ПОСТОЯННЫМ ПАРАМЕТРОМ)

Правило выполнения шага совпадает с (6.18), где $\delta_k \equiv \varepsilon$. Останов происходит в том случае, когда становится $W_{k+1}(y^{k+1}, \varepsilon) = 0$.

МЕТОД 2. (С НАСТРОЙКОЙ ПАРАМЕТРА δ_k)

Перед началом поиска для $k=0$ выбирается значение $\delta_k \gg \varepsilon$. Далее в процессе поиска при выполнении условия

$$W_{k+1}(y^{k+1}, \delta_k) > 0, \quad (6.19)$$

сохраняется прежнее значение параметра, т.е. полагается $\delta_{k+1} = \delta_k$. Если на некотором шаге $W_{k+1}(y^{k+1}, \delta_k) = 0$, то выполняется проверка величины δ_k . Если $\delta_k \leq \varepsilon$, то поиск останавливается, иначе полагают $\delta_k := \max\{\varepsilon, \delta_k/2\}$ и повторяют попытку выбора y^{k+1} из (6.18).

МЕТОД 3. (С ЭКСТРЕМАЛЬНЫМ ВЫБОРОМ δ_k)

На каждом шаге выбирают наибольшее значение δ_k^* , при котором для $\delta_k = \delta_k^* - 0$ еще выполняется (6.19), т.е.

$$\delta_k^* = \sup\{\delta_k : W_{k+1}(y^{k+1}, \delta_k) > 0\} \quad (6.20)$$

Останов происходит при $\delta_k^* \leq \varepsilon$.

Очевидно, что экстремальный принцип выбора параметра метода эквивалентен методу со следующим правилом выбора шага

$$y^{k+1} = \arg \max_{y \in Y} \rho(f_k^-(y), S(F_k)) \quad (6.21)$$

и с остановом при выполнении условия

$$\rho(f_k^-(y^{k+1}), S(F_k)) \leq \varepsilon \quad (6.22)$$



Замечание. Метод (6.21), (6.22) можно считать прямым обобщением метода Пиявского (см. раздел 4.3 четвертой главы) на многокритериальные липшицевы задачи.

Дополнительно отметим еще два момента, связанных с реализацией методов 1, 2, 3. Во-первых, в качестве оценки в целом всего множества эффективных точек Y° используется множество Y_k° , алгоритм получения которого, по результатам выполненных измерений, приведен перед разделом 6.1. Во-вторых, при построении оценок $f_k^\pm(y)$ вектор-функции $f(y)$ используются не точные значения констант липшица L_1, \dots, L_n , а их оценки, L_1^k, \dots, L_n^k полученные по результатам испытаний согласно правилам

$$\mu_i^k = \max_{\substack{j=1, \dots, k-1 \\ s=j+1, \dots, k}} \left(\frac{|f_i^j - f_i^s|}{\|y^j - y^s\|} \right)$$

$$L_i^k = \begin{cases} r \cdot \mu_i^k, & \text{if } \mu_i^k > 0 \\ 1, & \text{if } \mu_i^k = 0 \end{cases}, \quad r > 1.$$

Приведем без доказательства две теоремы о свойствах построенных методов. Доказательства можно найти в работе [18].

Теорема 6.4. Пусть Y —компакт, $f \in \Phi = Lip_L(Y)$, $\varepsilon > 0$, тогда в методах 1, 2, 3 за конечное число шагов выполнится критерий останова, и, если к этому моменту оценки констант Липшица окажутся не меньшими, чем их истинные значения, т.е. $L_i^k \geq L_i$ ($i=1, \dots, n$), построенное множество Y_k^o определит ε —оптимальное (по Слейтеру) решение задачи (6.1).

Теорема 6.5. Пусть Y —компакт, $f \in \Phi = Lip_L(Y)$ и точность $\varepsilon = 0$, тогда методы 2-й и 3-й порождают бесконечную последовательность измерений. Если при этом, начиная с некоторого шага, выполнится $L_i^k \geq L_i$ ($i=1, \dots, n$), то будут иметь место следующие свойства.

1. Для всякой эффективной точки $y^o \in Y^o$ найдется предельная точка y_∞^o последовательности оценок Y_k^o , что $f(y_\infty^o) \leq f(y^o)$.

2. Все предельные точки y_∞^o последовательности испытаний являются слабоэффективными, т.е. $\forall y_\infty^o: y_\infty^o \in Y^o$.

Если же, начиная с некоторого шага, $L_i^k > L_i$ ($i=1, \dots, n$), то множество $\{y_\infty^o\}$ предельных точек последовательности испытаний совпадают со множеством слабо эффективных точек, т.е. $\{y_\infty^o\} = Y^o$.

Заметим, что общая схема описанных в разделе 6.3 методов точно вычислительно не реализуема. Ее следует рассматривать как идеальную теоретическую схему. Вычислительно реализуемые версии этих методов могут быть получены с использованием компонентных подходов, описанных в конце главы пять, например, с использованием триангуляции области поиска по аналогии с методом [19], рассмотренным в пункте 5.7.4 пятой главы.

6.4. Метод построения равномерной оценки множества слабо эффективных точек

В этом разделе приводится описание и свойства алгоритма, предложенного в работе Маркина Д.Л., Стронгина Р.Г. [30] для решения скалярной минимаксной задачи (6.10), порождаемой непараметрической сверткой (6.8). Метод основан на редукции размерности в задаче (6.1) с использованием однозначных непрерывных отображений отрезка $[a, b]$ на гиперпараллелепипед $D \subset \mathbb{R}^N$, порождающих кривые или развертки Пеано. Эти отображения были рассмотрены в разделе 5.2. С их использованием задача (6.1) с функциями $f_i(y)$, удовлетворяющими условию Липшица с константами L_i , трансформируется в задачу

$$f(x) = (f_1(x), \dots, f_n(x)) \rightarrow \min, \quad x \in [a, b] \quad (6.23)$$

где соответствующие функции одного переменного $f_i(x)$ удовлетворяют на $[a, b]$ равномерному условию Гельдера (смотрите пункт 1.4.2.1. первой главы) вида

$$\forall x', x'' \in [a, b] \quad |f_i(x') - f_i(x'')| \leq K_i (|x' - x''|)^{1/N}, \quad (6.24)$$

где $K_i = 4L_i N^{1/2}$ ($i=1, \dots, n$).

Перепишем (6.8), (6.9) в несколько иной форме, а именно, введем функции

$$H(x, \tilde{x}) = \min \left\{ \left(f_s(x) - f_s(\tilde{x}) \right) / K_s : s = 1, \dots, n \right\}, \quad (6.25)$$

$$\varphi(x) = \max \{ H(x, \tilde{x}) : \tilde{x} \in [a, b] \}. \quad (6.26)$$

Заметим, что функция непараметрической свертки $\varphi(x)$ для задачи (6.23) соответствует ранее введенной в (6.8) функции $\psi(f(y))$ для задачи в форме (6.1).

На основании теоремы 6.2 приходим к скалярной задаче

$$\varphi^* = \min\{\varphi(x) : x \in [a, b]\}, \quad (6.27)$$

множество глобальных минимумов которой совпадает с множеством слабо эффективных точек задачи (6.23).

Приведем описание алгоритма решения задачи (6.27), следуя [30].

ОПИСАНИЕ АЛГОРИТМА

Первые две итерации осуществляются в концевых точках $x^0=a$, $x^1=b$ интервала $[a, b]$. Выбор каждой очередной точки x^{k+1} ($k > 1$) осуществляется по следующим правилам.

ПРАВИЛО 1. Точки выполненных итераций перенумеровываются нижним индексом в порядке возрастания координаты: $a=x_0 < x_1 < \dots < x_k=b$.

ПРАВИЛО 2. Вычисляются нижние оценки μ_s^k коэффициентов Гельдера для частных критериев $f_s(x)$ ($s=1, \dots, n$)

$$\mu_s^k = \max\{|f_s(x_i) - f_s(x_j)| / (x_i - x_j)^{1/N} : 0 \leq j < i \leq k\}.$$

Если μ_s^k оказывается равным нулю, полагается $\mu_s^k=1$.

ПРАВИЛО 3. Каждой точке x_i сопоставляется значение

$$z_i = \max\{H(x_i, x_j) : 0 \leq j \leq k\},$$

где $H(x_i, x_j)$ определяется формулой (6.25) с заменой K_s на оценки μ_s^k .

ПРАВИЛО 4. Для каждого интервала (x_{i-1}, x_i) ($1 \leq i \leq k$) вычислить характеристику

$$W(i) = \Delta_i + (z_i - z_{i-1})^2 / \Delta_i - 2(z_i + z_{i-1})/r,$$

где $\Delta_i = (x_i - x_{i-1})^{1/N}$, а $r > 1$ — параметр метода.

ПРАВИЛО 5. Выбрать первый из интервалов с наибольшей характеристикой

$$W(t) = \max\{W(i) : 1 \leq i \leq k\},$$

выполнить очередное измерение в точке

$$x^{k+1} = (x_t + x_{t-1})/2 - |z_t - z_{t-1}|^N \text{sign}(z_t - z_{t-1}) / (2r)$$

и, положив $k=k+1$, вернуться к правилу 1.

ПРАВИЛО 6. Итерации прекращаются при выполнении условия $\Delta_i \leq \varepsilon$.

Прежде, чем сформулировать теорему о свойствах построенной вычислительной процедуры, дадим определение.

Определение 6.2. Говорят, что последовательность $\{x^k\}$ равномерно сходится ко множеству слабо эффективных точек X^0 задачи (6.23), если $\lim_{k \rightarrow \infty} \rho(X^0, \{x^0, \dots, x^k\}) = 0$, где $\rho(X^0, \{x^0, \dots, x^k\}) = \sup_{x \in X^0} \inf_{0 \leq i \leq k} |x - x^i|$

Теорема 6.6. Пусть в задаче (6.23) функции $f_i(x)$ удовлетворяют условиям Гельдера (6.24) и на последовательности $\{x^k\}$, порожденной описанным алгоритмом (при $\varepsilon=0$), начиная с некоторого шага выполняются неравенства $r\mu_s^k > 4K_s$ ($s=1, \dots, n$). Тогда множество предельных точек последовательности $\{x^k\}$ совпадает с множеством слабо эффективных точек задачи, т.е. $\{x^\infty\} = X^0$. При этом сходимость $\{x^k\}$ к X^0 является равномерной.

ДОКАЗАТЕЛЬСТВО приведено в [30].

ПРИМЕР. В качестве иллюстрации на рис. 6.4 приведены точки 200 итераций, выполненных описанным алгоритмом для двумерной задачи ($N=2$) с двумя критериями ($n=2$) следующего вида

$$f_1(y_1, y_2) = \min\{(y_1^2 + y_2^2)^{1/2} + 0,5; ((y_1 - 1,5)^2 + (y_2 + 1,5)^2)^{1/2}\}$$

$$f_2(y_1, y_2) = ((x_1 + 0,5)^2 + (x_2 - 0,5)^2)^{1/2}.$$

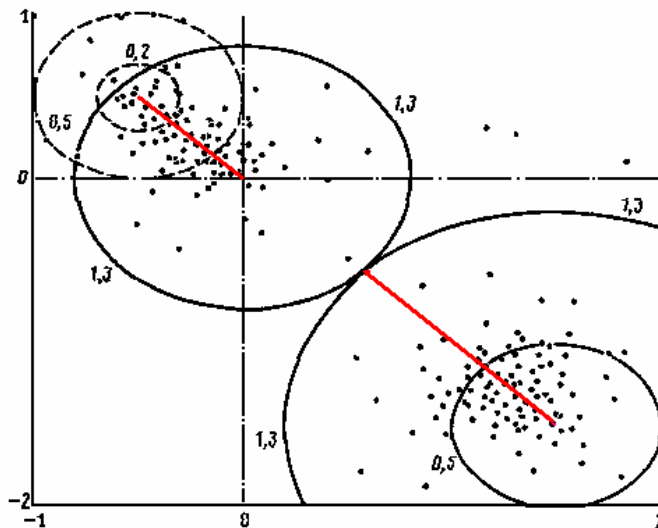


Рис. 6. 4. Размещение испытаний методом Стронгина Р.Г., Маркина Д.Л. в задаче с множеством X^0 в виде двух выделенных отрезков

На рисунке достаточно ясно проявляется сгущение точек испытаний в окрестности множества решений многокритериальной задачи.

Лист регистрации изменений

Дата	Автор	Комментарии
04.08.02	Городецкий С.Ю.	Создание документа
10.08.02	Городецкий С.Ю.	Внесение изменений
12.08.02	Городецкий С.Ю.	Внесение изменений
14.08.02	Городецкий С.Ю.	Включение рисунков
18.08.02	Городецкий С.Ю.	Внесение изменений
25.09.03	Городецкий С.Ю.	Исправление ошибок, изменение рисунка
12. 10.03	Городецкий С.Ю.	Окончательная редакция

Глава 7. Модели и методы поиска локально-оптимальных решений

В этом разделе лекционного курса мы вернемся к изучению следующей однокритериальной задачи математического программирования

$$f(y) \rightarrow \min, \quad y \in Y, \quad f: Y \rightarrow R^1, \quad (7.1)$$

$$Y = \{y \in D \subseteq R^N; g(y) \leq 0, h(y) = 0\}, \quad (7.2)$$

$$D = \{y \in R^N : a \leq y \leq b\}, \quad (7.3)$$

В некоторых случаях мы будем предполагать, что в определении допустимой области (7.2) ограничения–равенства отсутствуют, т.е. множество Y имеет вид

$$Y = \{y \in D \subseteq R^N; g(y) \leq 0\}. \quad (7.2')$$

Формально, как это уже отмечалось в четвертой главе, область (7.2) всегда может быть сведена к форме (7.2') за счет замены одного равенства $h(y)=0$ двумя неравенствами $h(y) \leq 0, -h(y) \leq 0$. Заметим, что при численном решении обычно вводится допустимая невязка $\delta > 0$ для равенств, и они заменяются совокупностью следующих условий: $h(y) - \delta \leq 0, -h(y) - \delta \leq 0$.

В разделе 1.2 первой главы были рассмотрены различные понятия решения, используемые для задач (7.1)–(7.3) или (7.1), (7.2'), (7.3). Последняя постановка (E – постановка) определяла задачу поиска локального экстремума.

Определение 7.1. *Задача поиска локального минимума (локально-оптимального решения) состоит в том, чтобы для заданной начальной точки $y \in Y$ найти точку и значение локального минимума $y^o(y)$, в области притяжения которого находится заданная начальная точка y .*

Существует достаточно много случаев, требующих поиска локально-оптимальных решений. Такая задача возникает, как правило, тогда, когда известна приближенная оценка y^o глобально-оптимального решения, найденная с неудовлетворительной точностью. В этом случае достаточно найти с высокой точностью локально-оптимальное решение, соответствующее начальной точке поиска y^o . Если эта точка была выбрана правильно, то найденный локальный минимум $y^o(y^o)$ зависящий от y^o , будет являться глобальным минимумом задачи.


В общем случае, вычислительные затраты на локальное уточнение решения оказываются меньшими тех затрат, которые были бы необходимы процедурам многоэкстремальной оптимизации для достижения той же точности, которую может дать локальный метод.

Следует указать и на другие ситуации, в которых целесообразно ставить задачу о поиске локального решения. Они могут возникнуть при необходимости предварительного исследования структуры решаемой задачи. Например, если из нескольких начальных точек метод локальной оптимизации получил существенно разные решения, то это говорит о многоэкстремальном характере задачи оптимального выбора и о необходимости применения к ее решению методов глобального поиска.

Наконец, следует учитывать, что в задачах высокой размерности регулярные методы поиска глобального решения не могут быть применены из-за чрезвычайно больших вычислительных затрат на покрытие области точками испытаний. Это остается справедливым даже в случае использования эффективных методов, строящих адаптивные существенно неравномерные покрытия области (такие методы были рассмотрены в пятой главе). В задачах высокой размерности практически единственным средством их решения остаются методы локальной оптимизации, совмещенные, как это было описано в разделе 5.1, с процедурами предварительного отбора начальных точек, используемых при локальной оптимизации.

Кроме перечисленных случаев методы локального поиска необходимы в задачах слежения за дрейфом выделенного экстремума при постепенном изменении структуры задачи (за счет влияния переменных параметров). Такая ситуация (см. пример E в разделе 1.2 первой главы) возникает, например, в методах внешнего штрафа, рассмотренных в разделе 3.2.

Еще один важный случай связан с решением одноэкстремальных задач, относящихся к специальным классам, допускающим применение принципов оптимальности при построении методов оптимизации. Наиболее важный — класс выпуклых задач. Для него возможна разработка специальных эффективных методов, как это будет показано в разделе 7.1.

 **Замечание.** При построении методов поиска локально-оптимальных решений всегда неявно предполагается, что задача является одноэкстремальной (по крайней мере в той подобласти, где применяется локальный метод, это действительно так). Поэтому в дальнейшем не будем проводить различий между локальным y^0 и глобальным y^* минимумами. Всюду в этом разделе будет использоваться обозначение y^* .

7.1. Применение принципов оптимальности при построении методов локальной оптимизации выпуклых гладких задач

Концепция применения принципов оптимальности (одношаговой оптимальности) к построению вычислительных процедур поиска минимума была рассмотрена в разделе 4.2 четвертой главы. Использование этих принципов возможно лишь в том случае, когда класс задач Φ (т.е. принятая модель задачи оптимального выбора — см. определение 1.7) допускает введение функции эффективности $V_k(\omega_k(f, g, h))$ для текущей оценки решения e^k , построенной по априорной и накопленной поисковой информации (см. разделы 4.2, 4.3, а также 1.4). Иначе говоря, нужно, чтобы по результатам ω_k конечного числа испытаний было возможно построение оценок положения решения в той или иной форме.

Для большинства классов задач локальной оптимизации, интересных с точки зрения практики, этого сделать нельзя, поскольку имеющаяся о них априорная информация бывает достаточно скудной (как правило — гладкость определенного порядка и некоторые дополнительные неформальные предположения). Этой информации обычно недостаточно для построения необходимых оценок. Такие задачи можно назвать *задачами локальной оптимизации общего вида*. Методы их решения будут рассмотрены в разделах 7.2 – 7.7.

Кроме общих можно выделить специальные классы задач, допускающие конечные оценки положения решения. Важным является класс выпуклых задач и его обобщения (см. раздел 1.4.1). Ниже будем рассматривать гладкие выпуклые задачи. Представленные здесь подходы можно обобщить на гладкие задачи с квазивыпуклыми и псевдовыпуклыми целевыми функциями.

Вопросы построения эффективных методов и оценки вычислительной трудоемкости для класса выпуклых гладких задач подробно рассмотрены в книге А.С. Немировского и Д.Б. Юдина [36]. В лекции этот материал рассматривается в кратком изложении.

Определение 7.2. Задача математического программирования (7.1)–(7.3) называется выпуклой, если f, g — непрерывные выпуклые функции на выпуклом D , а ограничения–равенства отсутствуют или афинны, т.е. имеют вид $h(y) = Ay + c$.

7.1.1. Метод центров тяжести

Далее будем предполагать, что ограничения–равенства отсутствуют, т.е. задача имеет форму (7.1), (7.2'), (7.3). Дополнительно потребуем, чтобы $f, g \in C^1(E)$ и было возможно проводить их испытания первого порядка, в результате которых в каждой точке $y^k \in D$ вычислялись значения этих функций $f^k = f(y^k), g^k = g(y^k)$ и градиенты $\nabla f^k = \nabla f(y^k), \nabla g^k = \nabla g(y^k)$.

Предположим, что введены нормировочные коэффициенты $c_0 > 0, c_1 > 0, \dots, c_m > 0$ для целевой функции f и компонент функции ограничений–неравенств $g = (g_1, \dots, g_m)$. Для выпуклой задачи любая локально–оптимальная точка будет одновременно являться ее глобально оптимальным решением y^* . Зададим относительную точность решения задачи $\varepsilon > 0$. Под решением с относительной точностью $\varepsilon > 0$ будем понимать любую точку $y^{*\varepsilon} \in D$, что

$$c_0 f(y^{*\varepsilon}) \leq c_0 f(y^*) + \varepsilon, \quad c_i g_i(y^{*\varepsilon}) \leq \varepsilon \quad (i=1, \dots, m).$$

В дальнейшем будем считать, что нормировка выполнена следующим образом

$$c_0 = (\sup\{f(y) : y \in D\} - \inf\{f(y) : y \in D\})^{-1},$$

$$c_i = (\max\{0; \sup\{g_i(y) : y \in D\}\})^{-1}, \quad (i=1, \dots, m).$$

Построим правило сокращения области поиска по результатам очередного k -го испытания. Обозначим начальную область размещения решения через $D_0 = D$, а последующие ее оценки — через D_k . Эти области будут играть роль текущих оценок решения. Такие оценки в разделе 4.2 были обозначены через e^k . Перед k -м испытанием имеем область D_{k-1} . Пусть в некоторой точке $y^k = w$ получены результаты k -го испытания. Возможно два случая, которые мы рассмотрим отдельно.

Случай А. $\exists i$, что $c_i g_i(y^k) \geq \varepsilon$.

По свойству 1.7 выпуклых функций (формула (1.34)) можно построить полупространство H_k^- вне которого не могут находиться ε -приближенные решения задачи

$$H_k^- = \{y \in R^N : c_i g_i(y^k) + c_i (\nabla g_i(y^k), y - y^k) \leq \varepsilon\}.$$

В этом случае область возможного размещения решения сокращается до пересечения

$$D_k = D_{k-1} \cap H_k^-.$$

Введем в качестве значения функции эффективности на шаге меру области D_k , понимаемую в смысле лебегова объема $V(D_k)$, т.е.

$$V_k(\omega_k(f, g)) = V(D_k).$$

Тогда при фиксированном направлении вектора нормали d к границе ∂H_k^- области отсекающая (т.е. при заданном векторе $\nabla g_i(w) = d$) наихудшее значение функции эффективности будет достигаться при $c_i g_i(w) = \varepsilon$, что соответствует прохождению отсекающей гиперплоскости ∂H_k^- точно через точку $y^k = w$ проведенного испытания.

Случай В. $\forall i: c_i g_i(y^k) < \varepsilon$, т.е. точка измерения допустима в пределах заданной погрешности. При этом отсечение части области поиска можно выполнить, используя результат испытания целевой функции. А именно, полупространство, вне которого не содержится ε -приближенных решений, будет иметь вид

$$H_k^- = \{y \in R^N: (\nabla f(y^k), y - y^k) \leq 0\}.$$

Новая оценка области возможного размещения решения, как и в случае А, примет вид

$$D_k = D_{k-1} \cap H_k^-.$$

В этом случае отсекающая гиперплоскость ∂H_k^- с нормалью $d = \nabla f(w)$ будет всегда проходить через точку испытания y^k (см. рис. 1.16). Таким образом, мы рассмотрели оба случая А и В.

Можно сделать следующий вывод: в общем случае, если результат испытания еще не известен, в наихудшем случае выполняется отсечение части области D_{k-1} гиперплоскостью $\partial H_k^-(w, d)$, проходящей через точку планируемого испытания w с неизвестной ориентацией нормали d . Гарантированное значение функции эффективности на шаге, согласно (4.31), соответствует значению точной верхней грани меры новой области, где верхняя грань берется по d .

$$W_k(w) = W_k(w, D_{k-1}) = \sup \{V(D_{k-1}) \cap H_k^-(w, d): \|d\| = 1\}.$$

Выбор точки очередного измерения по принципу наилучшего гарантированного на шаге результата определяется соотношением

$$y^k = \arg \min \{W_k(w, D_{k-1}): w \in D_{k-1}\}.$$

Оно определяет выбор точки испытания, приводящий (в наихудшем случае) к наиболее быстрому сокращению лебегова объема области возможного размещения решения.

Ясно, что y^k должна выбираться в «центре» выпуклой области D_{k-1} . Однако, точный выбор затруднителен. В качестве естественного кандидата на роль «центра» области D_{k-1} подходит ее центр тяжести. Поэтому будем выбирать

$$y^k = \frac{\int_{D_{k-1}} y \, dy}{\int_{D_{k-1}} dy}$$

Метод, использующий приведенное правило выбора точки испытания в сочетании с описанными выше правилами отсечения называется *методом центров тяжести* (МЦТ). Останов в методе выполняется в том случае, когда

$$V(D_k) < \varepsilon^N V(D),$$

где N — размерность пространства поиска.

Остановимся на некоторых свойствах метода. Известна верхняя оценка коэффициента сжатия лебегова объема выпуклой замкнутой области при отсечении через центр тяжести.

Лемма 7.1.¹ Если D — выпуклая замкнутая ограниченная область в R^N , а \bar{D} — его часть, полученная отсечением гиперплоскостью, проходящей через центр тяжести D , то для отношения лебеговых объемов справедлива следующая оценка.

$$V(\bar{D})/V(D) \leq v(N) = (1 - (N/(N+1))^N) \leq (e-1)/e,$$

где e — основание натурального логарифма.

Лемма показывает, что метод центров тяжести сходится со скоростью геометрической прогрессии с коэффициентом не хуже, чем $v(N) \leq 0,514$. Существуют более точные оценки, учитывающие заданную относительную точность решения ε . А.С. Немировским и Д.Б. Юдиным доказано следующее свойство [36].

Свойство 7.1. Для любой гладкой выпуклой задачи (7.1), (7.2), (7.3) метод МЦТ, настроенный на относительную точность ε ($0 < \varepsilon < 1$), требует не более чем

$$K^*(N, \varepsilon) = \lceil N \ln(\varepsilon^{-1}) / \ln(v(N)^{-1}) \rceil + 2$$

испытаний, где $\lceil x \rceil$ — наименьшее целое, большее или равное x .

7.1.2. Метод эллипсоидов

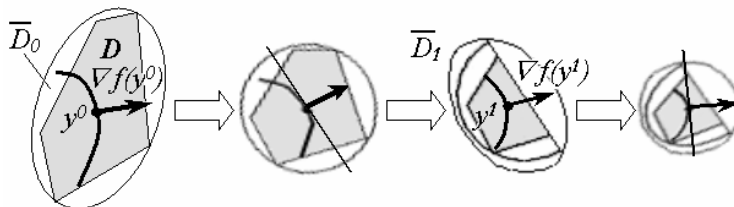
Метод центров тяжести вычислительно трудоемок из-за высоких вычисления затрат, связанных с вычислением центров тяжести сложных выпуклых областей D_{k-1} . А.С. Немировским и Д.Б. Юдиным предложена упрощенная вычислительно реализуемая модификация МЦТ, которую можно назвать методом эллипсоидов. Дадим ее концептуальное описание.

Пусть \bar{D}_0 — известный эллипсоид, содержащий начальную область поиска D . Идея метода заключается в том, чтобы на каждом шаге поиска область, полученную в результате очередного отсечения, заменять на содержащий ее эллипсоид \bar{D}_{k-1} минимального объема. Очередное измерение проводится в его легко вычисляемом геометрическом центре.

При этом, по сравнению с исходным методом МЦТ, требуется дополнительная проверка принадлежности точки y^k области D . Если $y^k \in D$, то отсекающая гиперплоскость строится обычным образом. Если же $y^k \notin D$, то в качестве ∂H_k выбирается любая гиперплоскость, опорная ко множеству точек

$$\{y: \rho(y, D) \leq \rho(y^k, D)\}.$$

Для упрощения алгоритма перед отсечением выполняется аффинное преобразование пространства, переводящее эллипсоид \bar{D}_{k-1} в шар.



¹ Б.С. Митягин Два неравенства для объемов выпуклых тел. – Мат. заметки, 1969, т.5, вып.1.

Таким образом, в новых координатах после отсечения всегда можно оставить область в виде полушария, которая заменяется на новый, описанный вокруг нее, эллипсоид минимального объема \bar{D}_k . Далее процесс повторяется.


Остановимся на свойствах. Известна простая геометрическая оценка коэффициента $(\beta(N))^N$, учитывающего дополнительное приращение объема за счет замены полушария описанным вокруг него эллипсоидом. Очевидно, что при $N=1$ $\beta(N)=1$. Оказывается, что при $N > 1$

$$\beta(N) = 2^{1/N} \left(N((N-1)/(N+1))^{1/2N} / (N^2 - 1)^{1/2} \right).$$

Свойство 7.2. Пусть $V(\bar{D}_0) / V(D) \leq \beta^*$ ($\beta^* > 1$), тогда для решения с относительной точностью ε любой гладкой выпуклой задачи методом эллипсоидов требуется не более чем

$$K^*(N, \beta, \varepsilon) = \lceil \ln(\beta^* / \varepsilon) / \ln 1/\beta(N) \rceil + 2$$

испытаний.

 **Замечание.** Метод эллипсоидов можно рассматривать как конкретный вариант алгоритма из класса методов растяжения пространства (R – алгоритмов) Н.З.Шора [49], рассматриваемых в разделе 7.5.

7.2. Принципы построения методов локальной оптимизации в задачах общего вида

В этом и следующих разделах изучаются методы поиска локально-оптимальных решений для классов задач, в которых невозможно применение принципа отсечений для сокращения (по проведенным измерениям) области поиска. Методы, используемые в таких задачах, можно отнести к траекторному типу. Они реализуют общую идею локального спуска и, в силу этого, имеют похожую структуру, которая рассматривается ниже.

Первоначально, в разделах 7.3 – 7.7, будут построены методы локальной оптимизации для простейшей задачи без каких-либо ограничений

$$f(y) \rightarrow \min, y \in R^N. \quad (7.4)$$

В них кроме целевой функции $f(y)$, определенной для $y \in R^N$, будем также задавать начальную точку поиска y^0 , необходимую для локальных методов. Несмотря на неограниченность области поиска всегда будем предполагать, что минимум в задаче существует.

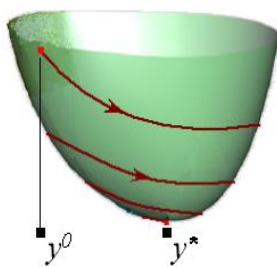
Позднее, в разделе 7.8 будут описаны специальные методы учета линейных ограничений, а также наиболее простых, но очень важных ограничений на переменные параллелепипедного типа $y \in D$.

Напомним, что общие способы учета произвольных функциональных ограничений были подробно изучены в разделе 3 части 2.

7.2.1. Общая структура методов поиска локального минимума, принцип локального спуска

К настоящему времени разработано огромное количество разнообразных методов локальной оптимизации для задач общего вида. Большинство из них используют *принцип локального спуска*, когда метод последовательно на каждом шаге переходит к точкам с меньшими значениями целевой функции. Почти все эти методы могут быть представлены в виде итерационного соотношения

$$y^{k+1} = y^k + x^k d^k, \quad (7.5)$$



где y^k — точки *основных испытаний*, состоящих в вычислении $I^k=I(y^k)$ — набора тех или иных *локальных характеристик* целевой функции в точке y^k , d^k — направления смещения из точек y^k , вычисляемые по результатам *основных испытаний*, а x^k — коэффициенты, определяющие величины смещений вдоль выбранных направлений.

В набор вычисляемых для функции локальных характеристик $I^k=I(y^k)$ могут входить: значение функции $f^k=f(y^k)$, вектор градиента $\nabla f^k=\nabla f(y^k)$, матрица вторых производных (гессиян) $\Gamma_k=\Gamma^f(y^k)$. Какой именно набор характеристик измеряется — зависит как от свойств решаемой задачи, так и от выбранного метода оптимизации.

Для определения величин смещений x^k вдоль направлений d^k методы могут выполнять вспомогательные (*рабочие*) шаги. Это приводит к дополнительным измерениям локальных характеристик целевой функции вдоль направления d^k (рис.7.1).

Переходы от точек y^k к точкам y^{k+1} выполняются таким образом, чтобы обеспечить существенное убывание значений функции $f^k=f(y^k)$ в результате шага. Заметим, что простого выполнения условия убывания значений f^k , когда для всякого k выполняется $f^{k+1} < f^k$, для обеспечения сходимости к решению задачи недостаточно.

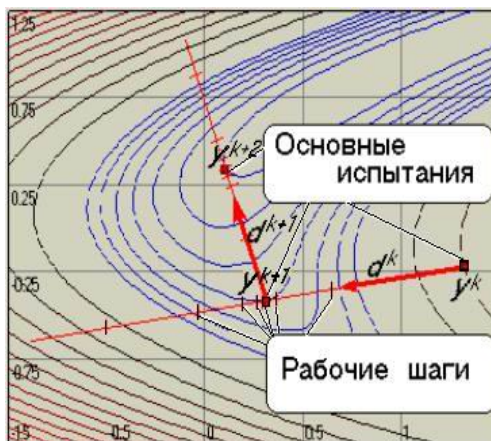


Рис. 7. 1. Общая структура организации поиска локального минимума

Останов вычислений в методах локальной оптимизации, применяемых в задачах без ограничений (7.4) с непрерывно дифференцируемыми целевыми функциями, происходит при выполнении условия достаточной малости нормы градиента

$$\|\nabla f(y^k)\| \leq \varepsilon. \quad (7.6)$$

Нужно отметить, что выполнение условия (7.6), в общем случае, не гарантирует близость точки y^k к решению задачи. Для методов, не использующих вычисление градиента (например, для методов прямого поиска Хука-Дживса или Нелдера-Мида), останов производится по другим правилам, своим для каждого из методов.

7.2.2. Измерения локальной информации и роль модели задачи в их интерпретации

Выбор направлений d^k при выполнении итераций методов локального поиска в большинстве методов происходит по результатам основных испытаний I^k . Для того, чтобы было возможно определить d^k , необходимо использовать имеющуюся априорную информацию или принятые предположения о свойствах решаемой задачи. Очевидно, что при отсутствии таких предположений, обоснованный выбор точек очередных испытаний был бы невозможен.

Совокупность предположений относительно свойств решаемой задачи мы назвали *моделью задачи* (определение 1.7 главы 1). В большинстве случаев, принятая модель задачи (7.4) может быть описана в терминах принадлежности целевой функции $f(y)$ некоторому классу функций: $f \in \Phi$.

Принятая модель задачи существенно влияет на интерпретацию результатов проведенных испытаний (см. примеры в разделе 1.4).

В задачах локальной оптимизации общего вида априорная информация о функциях бывает достаточно скудной. Например, если о функции f известно только то, что она непрерывно дифференцируема, т.е. принадлежит классу $\Phi = C(D)$, а испытание в точке y^k состоит в вычислении значения функции f^k и градиента ∇f^k , то нельзя указать никаких правил сокращения области поиска решения по результатам испытаний. Однако можно использовать информацию о векторе градиента для выбора направления d^k смещения текущей точки y^k . А именно, в качестве такого направления можно принять направление *антиградиента* $d^k = -\nabla f^k$, определяющего направление скорейшего локального убывания функции. Заметим, что если бы о целевой функции было известно больше, например, что $f \in C^2(D)$ и f близка к квадратичной функции, то при той же измеряемой информации можно было бы указать значительно лучшие стратегии выбора направлений поиска, чем антиградиент, используя квадратичную модель функции.

7.2.3. Классификация траекторных методов локального поиска

Обычно используют классификацию методов в зависимости от той локальной информации, которую метод получает при выполнении основных испытаний.

Определение 7.3. Если метод использует результаты испытаний, включающие вычисление производных функции до k -го порядка, то его относят к методам k -го порядка.

Обычно выделяют методы *второго порядка* (используют вычисления функции, ее градиента и матрицы Гессе), *первого порядка* (используют вычисления функции и ее градиента), а также *нулевого порядка* (используют только вычисления функции).

Определение 7.4. Если метод нулевого порядка не использует предположений о гладкости функции, то его называют методом прямого поиска.

Методы прямого поиска основаны на эвристических правилах определения направлений убывания минимизируемой функции и их структура может отличаться от описанной в пункте 7.2.1. Почти все остальные методы соответствуют структуре (7.5), и, следовательно, требуют для своей реализации

разработки специальных вычислительных процедур, позволяющих определять в (7.5) коэффициенты одномерных смещений x^k вдоль выбираемых направлений d^k .

7.2.4. Эффективные стратегии поиска вдоль направлений. Регуляризованные алгоритмы одномерного поиска

Вычислительная схема методов локального поиска (7.5) требует многократного применения процедур выбора одномерных смещений x^k вдоль направлений d^k .

В процессе работы метода эти процедуры могут выполняться сотни раз. Это накладывает повышенные требования к эффективности таких процедур. Остановимся на принципах и алгоритмах определения смещений.

Величина коэффициента x^k , определяющего длину шага вдоль направления d^k , может выбираться на основе различных критериев [2,4]. Весьма распространенными при построении методов являются следующие: критерий близости к минимуму по направлению, критерий существенности убывания функции, а также требование по степени уменьшения первоначального интервала возможных значений x^k .

В основе первого критерия лежит требование, чтобы в точке $y^k + x^k d^k$ величина скорости изменения функции f в направлении d^k была в заданное число раз меньше скорости ее изменения в точке y^k . Это требование формализуется следующим образом. Задается малый положительный коэффициент η , и величина x^k определяется условием

$$x^k \in \Pi_1(\eta), \quad 0 \leq \eta < 1 \quad (7.9)$$

$$\Pi_1(\eta) = \{x \geq 0 : |(\nabla f(y^k + x \cdot d^k), d^k)| \leq \eta \cdot (-\nabla f(y^k), d^k)\}. \quad (7.10)$$

В основе второго критерия (существенности убывания функции) лежит требование

$$x^k \in \Pi_2(\mu), \quad 0 < \mu < 1 \quad (7.11)$$

$$\Pi_2(\mu) = \{x \geq 0 : f(y^k + x \cdot d^k) \leq f(y^k) + \mu \cdot x \cdot (\nabla f(y^k), d^k)\}. \quad (7.12)$$

Эти два критерия используются совместно, и окончательное условие выбора x^k состоит в одновременном выполнении требований (7.9)-(7.12). При этом для их непротиворечивости на значения параметров μ и η накладывається дополнительное требование вида $\mu < \eta$, т.е.

$$x^k \in \Pi = \Pi_1(\eta) \cap \Pi_2(\mu), \quad 0 < \mu < \eta < 1. \quad (7.13)$$

Напомним, что используемые в формулах (7.10), (7.12) скалярные произведения градиента функции $f(y)$ на вектор направления d определяют производные функции f в точке y в направлении d , а именно $\partial f(y)/\partial d = (\nabla f(y), d)$.

На рис.7.2 дана иллюстрация выбора x^k на основе этих двух критериев. Если при выборе коэффициента одномерного шага используется условие (7.13) с $\eta \ll 1$, то говорят, что длина шага выбирается из условия "аккуратного" одномерного поиска.

При исследовании сходимости методов считается, что такой выбор соответствует выбору x^k из условия достижения минимума функции одномерного сечения

$$\varphi(x) = f(y^k + x d^k). \quad (7.14)$$

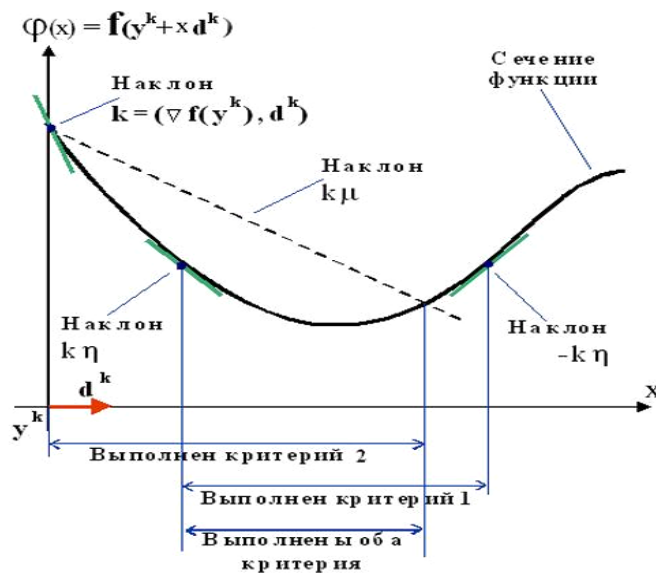


Рис. 7. 2 Критерии выбора коэффициента одномерного шага

Как организовать выбор x^k алгоритмически? Процесс выбора разбивается на два этапа. На первом этапе определяется промежуток $[0, X]$, на котором следует искать значение x^k . В задачах без явных ограничений на переменные этот промежуток должен иметь пересечение с искомым множеством Π из (7.13). Если же в задаче есть ограничения $y \in D$, то такого пересечения может не существовать, и тогда значение X определяется из условия попадания на границу области D из (7.3).

АЛГОРИТМ ОПРЕДЕЛЕНИЯ ПРОМЕЖУТКА $[0, X]$.

ШАГ 0. При решении задачи без ограничений (7.4) выбирается начальное значение $X^0 = \infty$. При решении задачи с ограничениями (7.1)–(7.3), включающими принадлежность точки y параллелепипеду D , по текущей точке y^k и направлению d^k начальное значение X^0 определяется как наименьшее значение x , при котором точка $y^k + x d^k$ попадает на границу этого параллелепипеда.

ШАГ 1. Выбирается малое $\delta > 0$ и $x = 0$.

ШАГ 2. Полагается $x = x + \delta$.

ШАГ 3. Если точка $y^k + x d^k \notin D$, то окончательно принимается $X = X^0$ и процесс останавливается. Если $y^k + x d^k \in D$ и $p = (\nabla f(y^k + x d^k), d^k) > 0$, т.е. обнаружено значение x , при котором функция f в одномерном сечении возрастает, то полагается $X = x$ и процесс останавливается. В противном случае удваивается величина шага $\delta := 2\delta$ и происходит возврат на шаг 2.

На втором этапе определения x^k на промежутке $[0, X]$ выполняется процедура поиска минимума функции одного переменного $\varphi(x)$ из (7.14), являющейся одномерным сечением функции $f(y)$ в направлении d^k . Поиск продолжается до момента первого попадания значения x в множество Π или же до того момента, когда будет достигнут заданный коэффициент сжатия σ ($0 < \sigma < 1$) для текущего интервала, содержащего решение, по отношению к длине исходного интервала $[0, X]$.

Дополнительное условие останова (по коэффициенту сжатия интервала) необходимо для задач с нарушением гладкости, а также для задач с ограничениями, в которых минимум может достигаться на границе интервала.

Заметим, что при поиске минимума функция $\varphi(x)$ всегда считается *унимодальной*, что позволяет применить известный *алгоритм золотого сечения*,

близкий к ε -оптимальному методу Фибоначчи ([46], а также [8] и [48]). Однако во многих задачах локальной оптимизации функции $\varphi(x)$ являются достаточно гладкими, что позволяет применять к поиску минимума алгоритмы, основанные на построении квадратичных аппроксимаций $\varphi(x)$ по результатам ее измерений. Такие методы называют *квазиньютоновыми*. Известно, что при определенных условиях они способны обеспечить скорость сходимости более высокого порядка, чем метод золотого сечения. Если же условия их сходимости будут нарушены, то такие методы могут расходиться.

Для осуществления одномерного поиска гладкой унимодальной функции наиболее подходящими являются *регуляризованные алгоритмы*, представляющие комбинацию метода золотого сечения с квазиньютоновым алгоритмом [4].

РЕГУЛЯРИЗОВАННАЯ ПРОЦЕДУРА ОДНОМЕРНОГО ПОИСКА состоит в следующем.

ШАГ 0. Полагаем $A=0, B=X, \tau = (-1+5^{1/2})/2, 0 < \delta < \tau$.

ШАГ 1. Выполняем три вычисления по методу золотого сечения:

ШАГ 1.1. Вычисляем $x_1 = B - (B-A)\tau, x_2 = A + \tau(B-A)$ и $\varphi_1 = \varphi(x_1), \varphi_2 = \varphi(x_2)$.

ШАГ 1.2. Если $\varphi_1 \leq \varphi_2$, то полагаем $B = x_2$ и $x_3 = B - (B-A)\tau, \varphi_3 = \varphi(x_3)$. При $\varphi_1 \leq \varphi_3$ полагаем $A = x_3, x = x_1$, иначе $B = x_1, x = x_3$. Если $\varphi_1 > \varphi_2$, то полагаем $A = x_1$ и $x_3 = A + (B-A)\tau, \varphi_3 = \varphi(x_3)$. При $\varphi_2 > \varphi_3$, полагаем $A = x_2, x = x_3$, иначе $B = x_3, x = x_2$. В результате выполнения шага 1 получаем три точки с вычисленными значениями функции, а также интервал $[A, B]$ с расположенной внутри него точкой x , соответствующей лучшему вычисленному значению функции.

ШАГ 2. Определяем u — точку измерения функции по квазиньютоновскому правилу. А именно, u определяется как точка минимума квадратичной аппроксимации функции $\varphi(x)$, построенной по значениям $x_1, \varphi_1; x_2, \varphi_2; x_3, \varphi_3$:

$$u = (-\varphi_1(x_3^2 - x_2^2) + \varphi_2(x_3^2 - x_1^2) - \varphi_3(x_2^2 - x_1^2)) / (2(\varphi_1(x_3 - x_2) - \varphi_2(x_3 - x_1) + \varphi_3(x_2 - x_1))).$$

ШАГ 3. Определяем v — точку измерения функции по правилу золотого сечения

$$v = \begin{cases} A + (x - A)\tau & \text{при } x > (A + B)/2 \\ B - (B - x)\tau & \text{при } x < (A + B)/2 \end{cases}.$$

ШАГ 4. Выбираем точку w для очередного вычисления функции: Если $u \in [\min(v; x); \max(v; x)]$; т.е. если точка квазиньютоновского шага u незначительно уклоняется от середины отрезка $[A, B]$, то в качестве точки нового измерения выбирается точка $w = u$, однако, чтобы предотвратить ее слишком близкое размещение к точке прежнего измерения x , ее положение корректируется

$$w = \begin{cases} u + \delta \cdot \text{sign}(u - x) & \text{при } |u - x| < \delta \\ u & \text{при } |u - x| > \delta \end{cases}.$$

Если же u не принадлежит указанному интервалу, полагаем $w = v$.

ШАГ 5. Вычисляем $\varphi_w = \varphi(w)$. Проверяем, принадлежит ли w множеству Π из (7.13). Если она принадлежит, или же достигнут предельный уровень сжатия интервала, т.е. $|B - A| < \sigma X$, то переходим на шаг 7, если нет, переходим на шаг 6.

ШАГ 6. Через y обозначим левую из точек w, x , (а через φ_y соответствующее ей значение функции), через z - правую из точек w, x (а через φ_z - соответствующее значение функции). Если $\varphi_y \leq \varphi_z$, то полагаем $B = z, x = y$. Если $\varphi_y > \varphi_z$, то полагаем $A = y, x = z$. Выделяем из точек x_1, x_2, x_3, w три точки с наименьшими значениями функции и обозначаем их через x_1, x_2, x_3 , а соответствующие им значения функции — через $\varphi_1, \varphi_2, \varphi_3$. Переходим на шаг 2.

ШАГ 7. Выполняем завершающие операции: вычисляем $\varphi(x)$ на концах интервала $[A, B]$, в качестве x^k выбираем ту из трех точек w, A, B , где достигается меньшее значение функции φ , останавливаем поиск.

7.3. Аппроксимационные принципы построения алгоритмов. Анализ свойств классического градиентного метода и метода Ньютона

Прежде чем перейти к изучению эффективных методов поиска локально-оптимальных решений, следует рассмотреть простейшие классические методы. К ним следует отнести градиентные методы и метод Ньютона [1,4–6,8,10].

Идея градиентного метода была высказана О. Коши в середине XVIII века, но еще в 40-ые годы XX века градиентные методы представлялись вполне достаточным средством практического решения задач.

Градиентные методы предельно просты. Их можно описать общим итерационным соотношением (7.5), имеющим вид $y^{k+1} = y^k + x^k d^k$, где направление смещения из точки y^k совпадает с направлением антиградиента $d^k = -\nabla f^k$.

В этом направлении дифференцируемая функция $f(y)$ локально (в бесконечно малой окрестности точки y^k) убывает быстрее всего, т.к. производная $\partial f(y^k)/\partial v^k$ функции f в точке y^k , вычисленная в некотором направлении v^k ($\|v^k\|=1$), может быть представлена в виде скалярного произведения $\partial f(y^k)/\partial v^k = (\nabla f(y^k), v^k)$, которое достигает своего минимума именно в направлении антиградиента.

Заметьте, что в общем случае направление антиградиента в точке y^k не совпадает с направлением на локальный минимум (рис.7.3). Более того, это направление не инвариантно по отношению к растяжениям пространства переменных.

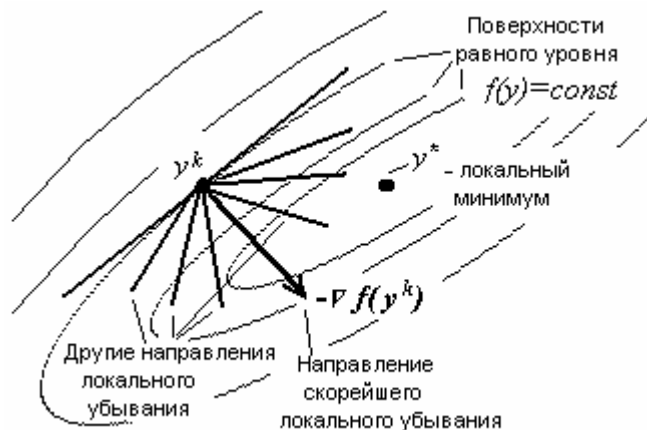


Рис. 7. 3. Отличие направления антиградиента от направления на локальный минимум

Существует несколько модификаций градиентного метода, различающихся правилом выбора величины смещения x^k в направлении антиградиента. Рассмотрим один из наиболее распространенных вариантов градиентного метода (метод *наискорейшего градиентного поиска*). Приведем правило выполнения шага в этом методе в том случае, когда в задаче оптимизации отсутствуют ограничения–неравенства из (7.2), то есть допустимая область $Y=D$ (общие принципы учета ограничений–неравенств рассмотрены в разделе 7.6). Применительно к этой задаче в методе наискорейшего градиентного поиска величина одномерного смещения x^k определяется из условия

$$f(y^k + x^k d^k) = \min \{f(y^k + x d^k) : x \geq 0, y^k + x d^k \in Y=D\}, \quad (7.15)$$

Таким образом, x — величина, определяющая смещение вдоль d^k , выбирается из условия достижения минимума функции f в области Y на луче $y^k + xd^k$, где $x \geq 0$. С вычислительной точки зрения правило (7.15) реализуется методом аккуратного одномерного поиска, описанным в пункте 7.2.4.

На градиентные методы можно также посмотреть с несколько иных позиций, а именно с позиций тех представлений о поведении минимизируемой функции, на которых основано правило выбора направления поиска. Градиентные методы основаны на локальной линейной модели функции $f(y)$ в окрестности точки y^k последнего испытания. Именно для линейной модели функции направление антиградиента является наилучшим с точки зрения задачи поиска минимума (для квадратичной модели это уже не так). Заметим, что методом используется не только локальная линейная модель. Выбор величины смещения по правилу (7.15) неявно предполагает нелинейность функции, особенно в задаче без ограничений.

В случае двух переменных правила поиска в методе наискорейшего градиентного поиска иллюстрирует рис.7.4. Следует заметить, что направления поиска на двух последовательных шагах d^k и d^{k+1} взаимно ортогональны, если решение вспомогательной задачи (7.15) достигается во внутренней точке допустимой области.

Сходимость процедур градиентного поиска может быть доказана при достаточно слабых предположениях о функции, минимум которой ищется. Одна из теорем о сходимости для случая использования метода при отсутствии ограничений (когда $D = R^N$) имеет следующую формулировку.

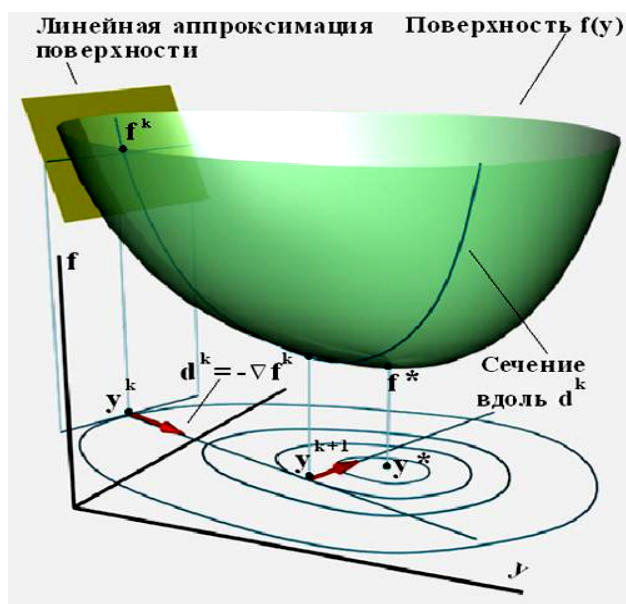


Рис. 7. 4. Наискорейший градиентный поиск

Теорема 7.1. Пусть в задаче (7.4) функция $f(y)$ непрерывно дифференцируема, ограничена снизу и ее градиент удовлетворяет условию Липшица с некоторой константой L , т.е. $\forall y', y'' \|\nabla f(y') - \nabla f(y'')\| \leq L \cdot \|y' - y''\|$. Тогда метод наискорейшего градиентного поиска для любой начальной точки y^0 строит последовательность y^k такую, что $\|\nabla f(y^k)\| \rightarrow 0$ при $k \rightarrow \infty$.

ДОКАЗАТЕЛЬСТВО [39] проведем следующим образом. Рассмотрим точки, лежащие на луче, порожденном направлением антиградиента $y(x) = y^k - x \nabla f(y^k)$. И сделаем верхнюю оценку для приращения функции, используя некоторую промежуточную точку $\theta \in [y(x) - y^k]$,

$$\begin{aligned} f(y(x)) - f(y^k) &= (\nabla f(\theta), y(x) - y^k) = (\nabla f(\theta) - \nabla f(y^k), y(x) - y^k) + \\ &+ (\nabla f(y^k), y(x) - y^k) \leq x^2 L \|\nabla f(y^k)\|^2 - x \|\nabla f(y^k)\|^2 = \varphi(x). \end{aligned}$$

$$\min_{x \geq 0} f(y(x)) - f(y^k) = f(y^{k+1}) - f(y^k) \leq \min_{x \geq 0} \varphi(x) = -\|\nabla f(y^k)\|^2 / (4L) < 0.$$

Таким образом, мы показали строго монотонное убывание последовательности значений функции на траектории поиска. В силу их ограниченности снизу, получаем, что $f(y^{k+1}) - f(y^k) \rightarrow 0$ при $k \rightarrow \infty$, но тогда из предыдущего неравенства $\|\nabla f(y^k)\| \rightarrow 0$ при $k \rightarrow \infty$. Теорема доказана.

При определенных дополнительных предположениях относительно функции f из утверждения теоремы будет следовать сходимость y^k к точке минимума y^* .

Сам факт сходимости еще не говорит об эффективности градиентного метода. Более того, ряд аналитических оценок показывают, что его эффективность в общем случае достаточно низка. Рассмотрим работу метода наискорейшего градиентного поиска на квадратичных функциях вида

$$f(y) = (y^T \Gamma y) / 2 + c^T y, \quad \Gamma^T = \Gamma \quad (7.16)$$

с положительно определенной матрицей Γ (т.е. строго выпуклых).

Известно, что на таких функциях метод наискорейшего градиентного поиска не обладает, в общем случае, конечной сходимостью (т.е. не определяет точку минимума за конечное число итераций), а порождает последовательность точек, сходящуюся к точке минимума y^* со скоростью геометрической прогрессии, знаменатель которой может быть близок к единице.

Теорема 7.2. Для квадратичной функции (7.16) с симметричной положительно определенной матрицей метод наискорейшего градиентного поиска сходится со скоростью геометрической прогрессии со знаменателем, не превосходящим значения q . При этом справедливы следующие оценки

$$\begin{aligned} \exists a = a(y^0), T > 0 : 0 \leq a \leq q &= (\lambda_{\min} / \lambda_{\max} - 1)^2 / (\lambda_{\min} / \lambda_{\max} + 1)^2, \\ f(y^k) - f(y^*) &\leq a^k (f(y^0) - f(y^*)), \\ \|y^k - y^*\| &\leq T a^{k/2} \|y^0 - y^*\|, \end{aligned}$$

где λ_{\min} и λ_{\max} — минимальное и максимальное собственные числа матрицы вторых производных $\Gamma^f = \Gamma$.

ДОКАЗАТЕЛЬСТВО теоремы можно найти в книге [10], здесь мы ограничимся его кратким изложением. Не снижая общности, будем считать, что в выражении для квадратичной функции (7.16) значение $c=0$.

Вначале получим оценку для q . Для этого достаточно исследовать один шаг алгоритма, поскольку остальные будут подобны первому. Пусть задана начальная точка y^0 . Найдем величину смещения x^1 на шаге. Введем обозначение $r = \nabla f(y^0)$. Из условия достижения минимума вдоль направления $-r$ имеем (при $c=0$)

$$\nabla f(y^0 - x^1 r) = \Gamma (y^0 - x^1 r) = r - x^1 \Gamma r = 0.$$

Умножая последнее равенство скалярно на r , выразим значение $x^1 = (r, r) / (r, \Gamma r)$. Отсюда уже несложно показать, что $f(y^1) = f(y^0) - 0,5(r, r)^2 / (r, \Gamma r)$, и, учитывая что при $c=0$ $f(y^*) = 0$, получить выражение

$$\begin{aligned} (f(y^0) - f(y^*)) / (f(y^0) - f(y^1)) &= 2(f(y^0) - f(y^*)) (r, \Gamma r) / (r, r)^2 = \\ &= ((r, \Gamma^{-1} r) (r, \Gamma r)) / (r, r)^2. \end{aligned}$$

Оценим его сверху, считая временно, что вектор r произволен. Представляя матрицу Γ в виде разложения $\Gamma = R^T D R$ с ортогональной матрицей R и диагональной матрицей D с элементами λ_i на диагонали. Выполняя замену $z = R r$, получим

$$((r, \Gamma^{-1} r) (r, \Gamma r)) / (r, r)^2 = ((z, D^{-1} z) (z, D z)) / (z, z)^2 = \Lambda(\alpha) M(\alpha),$$

где $\Lambda(\alpha) = \sum \lambda_i^{-1} \alpha_i$, $M(\alpha) = \sum \lambda_i \alpha_i$, $\alpha_i = z_i^2 / \sum_j z_j^2$.

Рассмотрим на плоскости переменных Λ, M область Q в виде выпуклой линейной оболочки семейства точек

$$Q = L((\lambda_{\min}; \lambda_{\min}^{-1}), \dots, (\lambda_{\max}; \lambda_{\max}^{-1})).$$

Найдем верхнюю оценку C произведения $\Lambda(\alpha) M(\alpha)$, аналитически решив экстремальную задачу

$$C = \max \{ \Lambda(\alpha) M(\alpha) : \alpha \geq 0, \|\alpha\| = 1 \} = \max \{ \Lambda M : (\Lambda; M) \in Q \}.$$

Из геометрического вида множества Q понятно, что решение находится в точке верхней части его границы. Записывая систему условий Куна–Таккера

$$\begin{cases} \Lambda = \mu, & M = \mu (1 / (\lambda_{\min} \lambda_{\max})) \\ M = 1 / \lambda_{\max} - (\Lambda - \lambda_{\max}) (1 / (\lambda_{\min} \lambda_{\max})) \end{cases}$$

находим, что $C = (\lambda_{\min} + \lambda_{\max})^2 / (4 \lambda_{\min} \lambda_{\max})$, следовательно

$$(f(y^0) - f(y^*)) / (f(y^0) - f(y^1)) \leq C$$

или

$$(f(y^1) - f(y^*)) \leq ((C - 1) / C) (f(y^0) - f(y^*))$$

Отсюда видим, что $q = ((C - 1) / C) = ((\lambda_{\min} / \lambda_{\max} - 1) / (\lambda_{\min} / \lambda_{\max} + 1))^2$. Требуемое выражение получено.

Чтобы получить оценку по координате полезно заметить, что квадратичная функция с $\lambda_{\min} > 0$ является сильно выпуклой (см. определение 1.10 пункта 1.4.1.1) со значением параметра сильной выпуклости $\rho = 0,5 \lambda_{\min}$ (по свойству 1.8'). Используя свойство сильной выпуклости (а именно — следствие 1.7'.2) и полученную оценку для ошибки по значению функции, приходим к неравенству

$$\|y^k - y^*\| \leq ((f(y^k) - f(y^*)) / \rho)^{1/2} \leq a^{k/2} ((f(y^0) - f(y^*)) / \rho)^{1/2}.$$

Представляя начальную разность $(f(y^0) - f(y^*))$ в виде разложения в окрестности y^* и оценивая сверху его квадратичный член, получим $(f(y^0) - f(y^*)) \leq \lambda_{\max} \|y^0 - y^*\|^2$. Эти две оценки вместе дают нужное неравенство для ошибки по координате при $T = (2 \lambda_{\max} / \lambda_{\min})^{1/2}$. Теорема доказана.

Из оценок теоремы 7.2 следует, что конечная сходимость из любой начальной точки y^0 возможна только при $q=0$ ($\lambda_{\min} = \lambda_{\max}$), когда поверхности равного уровня функции f являются сферами.

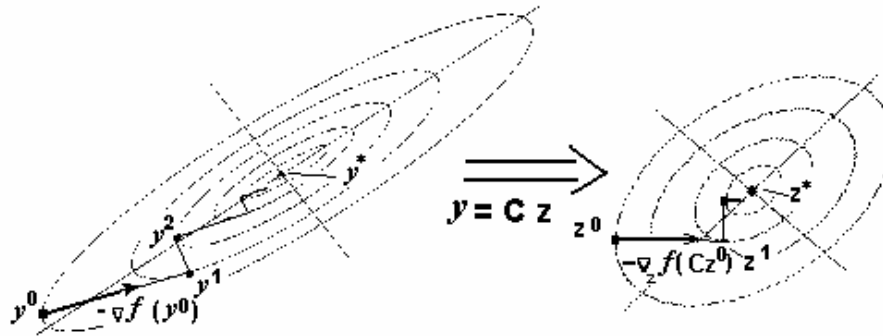


Рис. 7. 5. Влияние масштабирования на скорость сходимости наискорейшего градиентного поиска

Если же поверхности равного уровня сильно вытянуты (рис.7.5), что соответствует $\lambda_{\min} \ll \lambda_{\max}$, то значение q в оценке скорости сходимости будет близко к единице, и скорость сходимости к решению может оказаться чрезвычайно низкой, за исключением точек y^0 , лежащих на главных осях эллипсоидов $f(y)=const$ (в этих точках $a(y^0)=0$).

Таким образом, высокая скорость сходимости градиентного метода может быть обеспечена только за счет предварительного масштабирования задачи, т.е. выполнения такой замены переменных $y=Cz$, которая приводила бы (в новых переменных) к выполнению условия $\lambda_{\min} \approx \lambda_{\max}$.

Как уже отмечалось выше, направление антиградиента не инвариантно по отношению к линейным заменам переменных. Если выполнен переход к переменным $z: y=Cz$ (где C — матрица преобразования), то градиент функции в новых переменных z может быть вычислен как

$$\nabla_z f(Cz) = C^T \nabla f(y). \quad (7.17)$$

То же правило пересчета сохраняется, очевидно, и для антиградиента. Если перевести антиградиентное направление, вычисленное в пространстве z , в направление в пространстве старых переменных y , то получим скорректированное направление поиска

$$\bar{d}^k = C(-\nabla_z f(Cz^k)) = CC^T(-\nabla f(y^k)). \quad (7.18)$$

Возникающая матрица CC^T для коррекции направления антиградиента легко вычисляется для строго выпуклых квадратичных функций вида (7.16). Действительно, пусть в точке y^k для $f(y)$ измерено значение f^k , градиент ∇f^k и матрица вторых производных $\Gamma_k^f = \Gamma$. В силу равенства нулю всех производных выше второго порядка $f(y)$ из (7.16) совпадает со своей квадратичной аппроксимацией $P^k(y)$, построенной по измерениям f^k , ∇f^k , Γ_k^f , выполненным в точке y^k

$$P^k(y) = (y - y^k)^T \Gamma_k^f (y - y^k) / 2 + (\nabla f^k, (y - y^k)) + f^k. \quad (7.19)$$

Условие, определяющее y^* — точку минимума для $P^k(y)$, примет вид

$$\nabla P^k(y^*) = \Gamma_k^f (y^* - y^k) + \nabla f^k = 0, \quad (7.20)$$

откуда $y^* - y^k = (\Gamma_k^f)^{-1} (-\nabla f^k)$. Сравнивая полученное направление с направлением (7.18), видим, что в качестве матрицы преобразования в (7.18) можно использовать $CC^T = (\Gamma_k^f)^{-1}$.

Если бы $f(y)$ была произвольной дважды непрерывно дифференцируемой функцией, то в (7.19) квадратичная аппроксимация $P^k(y)$ уже не совпадала бы с исходной функцией, а условие (7.20) определяло бы лишь стационарную точку

для этой аппроксимации. Если именно в этой точке проводить очередное измерение локальных характеристик функции f (значения, градиента и матрицы вторых производных), приняв ее за y^{k+1} , получим классический метод Ньютона, имеющего вид итерационного соотношения

$$y^{k+1} = y^k + (\Gamma^f(y^k))^{-1} (-\nabla f(y^k)) \quad (7.21)$$

(направление шага $d^k = (\Gamma^f(y^k))^{-1} (-\nabla f(y^k))$, коэффициент длины шага $\alpha^k = 1$). Полезно обратить внимание на то, что этот метод, выполняя на каждом шаге некоторое преобразование пространства переменных. Построенное им направление d^k соответствует антиградиентному направлению функции f , если его вычислить в преобразованном пространстве.

На метод Ньютона можно посмотреть с другой точки зрения. А именно, правило итерации (7.21) основано на использовании квадратичной модели поведения функции $f(y)$, минимум которой ищется. Использование квадратичной модели приводит к отказу от использования направления антиградиента функции и применению вместо него скорректированного антиградиентного направления, приводящего в результате шага в стационарную точку текущей квадратичной аппроксимации функции.

Геометрическая интерпретация правила выбора направления поиска в методе Ньютона для выпуклой квадратичной функции приведена на рис.7.6.

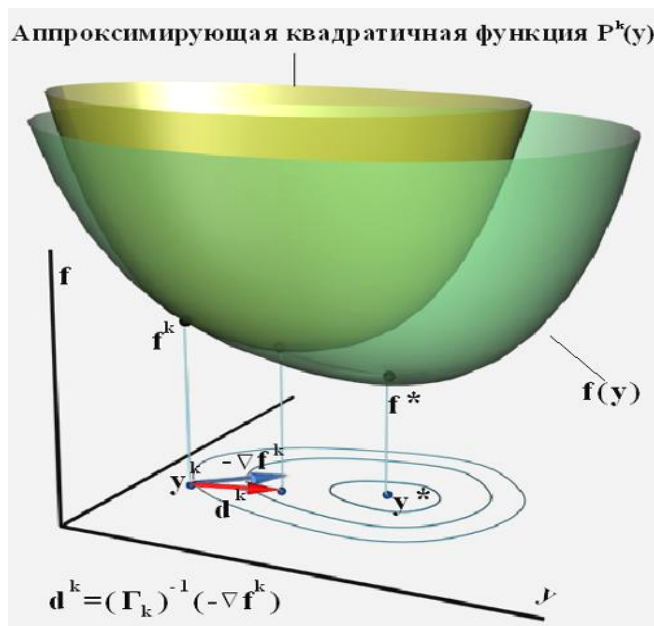


Рис.7. 6. Выбор направления в методе Ньютона

Что можно в достаточно общем случае сказать о сходимости метода Ньютона? Ответ дает следующая теорема.

Теорема 7.3. Для дважды непрерывно дифференцируемых функций (т.е. для $f \in C^2(D)$) с невырожденной матрицей $\Gamma^f(y^*)$ всегда существует такая ε — окрестность стационарной точки y^* функции $f(y)$, что для любой начальной точки y^0 из этой окрестности метод Ньютона будет сходиться сверхлинейно. Если при невырожденности матрицы $\Gamma^f(y^*)$ функция трижды непрерывно дифференцируема, то для начальных точек из некоторой ε — окрестности y^* метод сходится к y^* квадратично.

Вначале поясним терминологию.

Определение 7.5. *Линейной сходимостью называют сходимость по закону геометрической прогрессии.*

Линейная сходимость характерна для метода наискорейшего градиентного поиска (теорема 7.2).

Определение 7.6. *Говорят, что метод сходится сверхлинейно, если $\exists k > 0$ и последовательность чисел $\alpha_{k+1}, \dots, \alpha_{k+m}$ из интервала $(0, 1)$, стремящаяся к 0 при $m \rightarrow \infty$, что $\forall m > 0$ будет выполнено неравенство $\|y^{k+m} - y^*\| \leq \alpha_{k+1} \cdot \dots \cdot \alpha_{k+m} \|y^k - y^*\|$.*

Определение 7.7. *Говорят, что метод сходится квадратично, если $\exists T > 0$, что $\|y^{k+1} - y^*\| \leq T \|y^k - y^*\|^2$ при $\|y^0 - y^*\| < \varepsilon$.*

ДОКАЗАТЕЛЬСТВО теоремы 7.3. Предположим, что $f \in C^2(D)$. Тогда возможно следующее представление


$$\begin{aligned} \nabla f(y^k) &= 0 + \Gamma^f(y^*)(y^k - y^*) + o(\|y^k - y^*\|) \\ \Gamma^f(y^*) &= \Gamma^f(y^k) + \alpha(\|y^k - y^*\|), \end{aligned}$$

где $\alpha(\Delta)$ – бесконечно малая при $\Delta \rightarrow 0$. Подставляя второе выражение в первое и умножая на $(\Gamma^f(y^k))^{-1}$ получим

$$y^* - y^{k+1} = y^* - (y^k + (\Gamma^f(y^k))^{-1}(-\nabla f(y^k))) = o(\|y^k - y^*\|).$$

Следовательно, при достаточной близости начальной точки к стационарной найдется такая бесконечно малая $q = q(\|y^k - y^*\|)$, что $\|y^{k+1} - y^*\| \leq q \cdot \|y^k - y^*\|$. Это обеспечивает сверхлинейную сходимость.

В случае гладкости f до третьего порядка, можно гарантировать более высокий порядок малости у бесконечно малых в предыдущих выкладках, за счет чего в малой окрестности стационарной точки будет обеспечена квадратичная сходимость.

 **Замечание.** *Квадратично сходящаяся последовательность обладает скоростью сходимости более высокой (точнее говоря, более высокого порядка), чем у любой геометрической прогрессии. Сверхлинейная сходимость занимает промежуточное положение между квадратичной и линейной.*

Таким образом, из теоремы 7.3 следует, что при достаточно общих условиях метод Ньютона обладает тем, чего лишены градиентные методы — высокой скоростью сходимости. Однако это свойство сохраняется только в некоторой (заранее не известной!) окрестности решения. Вне этой окрестности метод Ньютона может вообще расходиться. Кроме того, итерация (7.21) требует обращения матрицы. Существенно также то, что в прикладных задачах достаточно часто встречаются ситуации, когда в точках последовательности y^k матрица $\Gamma^f(y^k)$ оказывается отрицательно определенной, знаконеопределенной или вырожденной. В последнем случае итерация (7.21) неприменима. Если же $\Gamma^f(y^*)$ не вырождена, но не знакоположительна, то, как следует из приведенной выше теоремы, метод Ньютона может сходиться к стационарной точке функции f ,

7.4. Эффективные методы второго порядка для гладких задач

В этом разделе изучаются методы локальной оптимизации, которые вычисляют в точке поиска y^k значения $f(y^k)$, $\nabla f(y^k)$, $\Gamma^f(y^k)$, т. е. явно используют значения матриц вторых производных. Эти методы являются улучшенными модификациями метода Ньютона, расширяющими область его сходимости и включающими методы борьбы с нарушением знакоположительности матриц Гессе.

7.4.1 Расширение области сходимости метода Ньютона за счет регулировки величины шага

Классический метод Ньютона обладает тремя существенными недостатками: возможной расходимостью для начальных точек, взятых вне некоторой окрестности решения, неприменимостью при вырождении матрицы вторых производных минимизируемой функции и возможной сходимостью к точкам максимумов или седловых точек в случае знакоотрицательности или знаконеопределенности этих матриц. Эти недостатки могут быть преодолены за счет модификаций метода Ньютона.

Первая модификация связана с изменением правила выбора длины шага. Ее изучению посвящен данный раздел.


В классическом методе Ньютона коэффициент длины шага $x^k \equiv 1$. В модифицированном методе (методе с регулировкой шага) x^k выбирается по алгоритму "аккуратного" одномерного поиска, описанному в пункте 7.2.4. Это приводит к сходимости из любой начальной точки для достаточно широкого класса функций и сохранению высокой (обычно сверхлинейной) скорости сходимости в окрестности решения. Метод Ньютона с регулировкой шага называют *методом Ньютона–Рафсона*.

Рассмотрим свойства данного метода. Выберем класс Φ одноэкстремальных тестовых функций, часто используемый для изучения свойств методов локального поиска. В качестве Φ возьмем класс $\Phi_{m,M}$ дважды непрерывно дифференцируемых функций, обладающих тем свойством, что $\exists m > 0$ и $M < \infty$, $m < M$, что $\forall y \in R^N, z \in R^N$

$$m\|z\|^2 \leq z^T \Gamma^f(y)z \leq M\|z\|^2. \quad (7.22)$$

Такие функции будут сильно выпуклы (см. свойство 1.8').

Можно показать, что условие (7.22) равносильно тому, что все собственные числа $\lambda_1(y), \dots, \lambda_N(y)$ матриц $\Gamma^f(y)$ лежат между m и M . Для этого достаточно заменить матрицу Гессе ее разложением через ортогональную R и диагональную матрицы (с собственными числами на диагонали), а затем перейти к новым переменным $w = Rz$.

 **Замечание.** Условие (7.22) гарантирует положительную определенность матрицы $\Gamma^f(y)$, достаточную для сходимости метода Ньютона к минимуму $f(y)$ из любой начальной точки, выбранной в достаточной близости от него. Однако из произвольно выбранной точки y^0 метод Ньютона для функции $f(y)$ из $\Phi_{m,M}$ может не сходиться.

ДОКАЗАТЕЛЬСТВО. Достаточно привести контр пример. Рассмотрим скалярный случай, когда $y \in R^1$. Построим четную функцию с минимумом в точке 0 (ее производная при этом будет функцией нечетной), для которой при $y^0 > 0$

следующая точка $y^1 = y^0 - f_0' / f_0''$ равнялась бы $(-y^0)$. В силу нечетности первой производной и четности второй производной, обязательно на следующем шаге выполнится равенство $y^2 = -y^1 = y^0$. Поэтому сходимости к точке минимума из точки y^0 не будет.

Таким образом, для классического метода Ньютона выбранный тестовый класс, с точки зрения сходимости из любой начальной точки, хорошим не является. Если бы удалось доказать, что метод Ньютона–Рафсона обладает на этом классе сходимостью из любой точки, это означало бы, что данный метод является улучшенной модификацией метода Ньютона.

Предварительно укажем на некоторые свойства выбранного класса функций.

Свойство 7.3. Любая функция $f(y)$ из класса $\Phi_{m,M}$ имеет единственный минимум.

Это непосредственно следует из сильной выпуклости функции (смотрите пункт 1.4.1.1).

Свойство 7.4. Для функций $f(y)$ из класса $\Phi_{m,M}$ существует взаимосвязь между ошибкой по координате и ошибкой по значению функции, выражаемая соотношением

$$0,5m\|y-y^*\|^2 \leq f(y) - f(y^*) \leq 0,5M\|y-y^*\|^2.$$

ДОКАЗАТЕЛЬСТВО сразу вытекает из (7.22) и разложения

$$f(y) - f(y^*) = (y-y^*)^T \Gamma^f(\theta) (y-y^*) / 2.$$

Свойство 7.5. Для функций $f(y)$ из класса $\Phi_{m,M}$ существует взаимосвязь между ошибкой по значению функции и нормой градиента

$$m(1+m/M)(f(y) - f(y^*)) \leq \|\nabla f(y)\|^2.$$

ДОКАЗАТЕЛЬСТВО. Раскладывая $f(y^*)$ в окрестности точки y , получим для приращения на классе $\Phi_{m,M}$

$$f(y) - f(y^*) = (\nabla f(y), y-y^*) - (y-y^*)^T \Gamma^f(\theta) (y-y^*) / 2 \leq \|\nabla f(y)\| \|y-y^*\| - 0,5m\|y-y^*\|^2.$$

Усиливая это неравенство слева с использованием свойства 7.4, придем к оценке

$$\|y-y^*\| \leq \|\nabla f(y)\| / m.$$

Используя еще раз два предыдущих неравенства, а также правое неравенство из свойства 7.4, получим

$$f(y) - f(y^*) \leq \|\nabla f(y)\|^2 / m - 0,5m\|y-y^*\|^2 \leq \|\nabla f(y)\|^2 / m - (m/M)(f(y) - f(y^*)).$$

Отсюда непосредственно вытекает справедливость свойства.

Поведение метода Ньютона с регулировкой шага на классе $\Phi_{m,M}$ определяется следующей теоремой [39]. Ее доказательство опирается на приведенные выше свойства.

Теорема 7.4 Метод Ньютона с регулировкой шага на функциях $f \in \Phi_{m,M}$ для любой начальной точки y^0 порождает последовательность точек y^k , сходящуюся к точке минимума y^* со сверхлинейной скоростью.

ДОКАЗАТЕЛЬСТВО распадается на две части. Во-первых, устанавливается сам факт сходимости (рассуждения здесь во многом аналогичны тем, которые использовались при обосновании сходимости градиентного метода). Во-вторых, оценивается скорость сходимости, исходя из того, что любая начальная точка гарантированно приводит в сколь угодно малую окрестность решения. Это

позволяет использовать в ней подходящие разложения, исходя из порядка гладкости функций. В приводимом доказательстве, прежде всего, представляет интерес техника подобного анализа.

Для доказательства сходимости рассмотрим точки, лежащие на луче, порожденном направлением Ньютона $d^k = (\Gamma^f(y^k))^{-1}(-\nabla f(y^k))$, а именно: $y(x) = y^k - x d^k$. Получим верхнюю оценку для приращения функции, используя некоторую промежуточную точку $\theta \in [y(x); y^k]$,

$$\begin{aligned} f(y(x)) - f(y^k) &= (\nabla f(y^k), y(x) - y^k) + 0,5 (y(x) - y^k)^T \Gamma^f(\theta) (y(x) - y^k) \leq \\ &\leq x(\nabla f(y^k), d^k) + x^2 0,5M \|d^k\|^2 = \varphi(x). \end{aligned}$$

$$\begin{aligned} f(y^{k+1}) - f(y^k) &= \min_{x \geq 0} f(y(x)) - f(y^k) \leq \min_{x \geq 0} \varphi(x) = \varphi\left(-(\nabla f(y^k), d^k) / (M \|d^k\|^2)\right) = \\ &= -(\nabla f(y^k), d^k)^2 / (2M \|d^k\|^2) = \\ &= -(\nabla f(y^k), (\Gamma^f(y^k))^{-1} \nabla f(y^k)) (\Gamma^f(y^k) d^k, d^k) / (2M \|d^k\|^2) \leq \\ &\leq -(m / (2M^2)) \|\nabla f(y^k)\|^2 \leq 0. \end{aligned}$$

Таким образом, либо норма градиента равна нулю (для функций класса $\Phi_{m,M}$ это означает точное определение глобального минимума), либо изменение функции на шаге строго отрицательно, т.е. имеет место строгое монотонное убывание последовательности значений функции на траектории поиска. В силу ограниченности снизу значений функции из класса $\Phi_{m,M}$, получаем, что $f(y^{k+1}) - f(y^k) \rightarrow 0$ при $k \rightarrow \infty$, но тогда из предыдущего неравенства $\|\nabla f(y^k)\| \rightarrow 0$ при $k \rightarrow \infty$. В рассматриваемом классе это означает сходимость по координате (см. доказательство свойства 7.5). Сходимость к решению из любой начальной точки доказана.

Оценим скорость сходимости. Будем использовать сокращенные обозначения: $f^k, \Gamma_k, \Gamma(\theta)$. Вначале рассмотрим точку, в которую перешел бы метод, выполнив шаг по методу Ньютона: $\tilde{y}^{k+1} = y^k - \Gamma_k^{-1} \nabla f^k$. Далее имеют место следующие оценки, первая из которых дважды использует свойство 7.4.

$$\|y^{k+1} - y^*\|^2 \leq 2(f^{k+1} - f^*)/m \leq 2(\tilde{f}^{k+1} - f^*)/m \leq (M/m) \|\tilde{y}^{k+1} - y^*\|^2.$$

$$\begin{aligned} \|\tilde{y}^{k+1} - y^*\|^2 &= (\tilde{y}^{k+1} - y^*, \tilde{y}^{k+1} - y^*) = (y^k - y^* - \Gamma_k^{-1}(\nabla f^k - \nabla f^*), \tilde{y}^{k+1} - y^*) = \\ &= ((E - \Gamma_k^{-1} \Gamma(\theta))(y^k - y^*), \tilde{y}^{k+1} - y^*) \leq \|\Gamma_k^{-1}\| \|\Gamma_k - \Gamma(\theta)\| \|y^k - y^*\| \|\tilde{y}^{k+1} - y^*\|. \end{aligned}$$

Или

$$\|y^{k+1} - y^*\| \leq \alpha_{k+1} \|y^k - y^*\|,$$

где $\alpha_{k+1} = (\sqrt{M/m/m}) \|\Gamma_k - \Gamma(\theta)\|$. Причем $\alpha_k \rightarrow 0$ при $k \rightarrow \infty$, тем самым выполнено определение 7.6 сверхлинейной сходимости. Теорема доказана.



Замечание. Если дополнительно потребовать от функции $f \in \Phi_{m,M}$ существования непрерывных третьих производных, можно доказать, что метод Ньютона–Рафсона будет сходиться квадратично.

Таким образом, модификация Ньютона–Рафсона расширяет область сходимости метода Ньютона.


7.4.2. Стратегии модификации матриц Гессе при нарушении их положительной определенности

Вторая модификация метода Ньютона связана с преодолением случаев отсутствия положительной определенности матрицы вторых производных. Следует обратить внимание на то, что при нарушении положительной определенности Γ_k (а значит и $(\Gamma_k)^{-1}$, если она существует) направление смещения $d^k = (\Gamma_k)^{-1}(-\nabla f(y^k))$ может не быть направлением убывания. Действительно, производная функции f в точке y^k по направлению d^k оценивается следующим образом

$$\partial f(y^k)/\partial d^k = (\nabla f(y^k), d^k) = -(\nabla f(y^k), (\Gamma_k)^{-1} \nabla f(y^k)).$$

В рассматриваемом случае знак этого произведения не определен и может оказаться положительным. В этом случае метод Ньютона–Рафсона применить нельзя, т.к. при его использовании смещение вдоль d^k будет нулевым.

ПРОСТОЙ АЛГОРИТМ. Простейший выход из этого положения может состоять в проверке знака скалярного произведения $(\nabla f(y^k), d^k)$ на каждом шаге. Если окажется, что $(\nabla f(y^k), d^k) < 0$, то выполняется обычный шаг по методу Ньютона–Рафсона, если же $(\nabla f(y^k), d^k) \geq 0$, то проводится замена направления поиска на антиградиентное направление $d^k = -\nabla f(y^k)$ и выполняется шаг по методу наискорейшего градиентного спуска.

 **Замечание.** Рассмотренная простая стратегия часто не оправдывает себя, поскольку при наличии овражности в области знаконеопределенности матриц Гессе простой алгоритм вырождается в метод наискорейшего градиентного поиска, очень плохо работающий в оврагах.

Таким образом, необходима более гибкая модификация метода. Необходимо как-то использовать матрицу Γ_k , хотя непосредственно ее применить нельзя. Основная идея, на основе которой выполняется модификация матриц, состоит в том, чтобы заменить Γ_k на достаточно близкую к ней (в смысле некоторой нормы) положительно определенную матрицу $\bar{\Gamma}_k$ и затем использовать ее в итерационном соотношении метода Ньютона–Рафсона

$$\begin{aligned} y^{k+1} &= y^k + x^k \bar{d}^k \\ \bar{d}^k &= (\bar{\Gamma}_k)^{-1}(-\nabla f(y^k)) \\ x^k &\in \Pi_1(\eta) \cap \Pi_2(\mu). \end{aligned} \quad (7.23)$$

Переход от Γ_k к положительно определенной матрице $\bar{\Gamma}_k$ обычно выполняется с помощью факторизации Γ_k , т.е. разложения ее в произведение матриц определенного вида.

Наиболее естественным представляется использование *спектрального разложения*. Определим для Γ_k набор собственных чисел $\lambda_1, \dots, \lambda_N$ и систему ортонормированных собственных векторов u^1, \dots, u^N . Тогда возможно следующее представление

$$\Gamma_k = \lambda_1 u^1 (u^1)^T + \dots + \lambda_N u^N (u^N)^T = \mathbf{U} \mathbf{L} \mathbf{U}^T,$$

где матрица \mathbf{U} составлена из вектор–столбцов u^1, \dots, u^N , \mathbf{L} – диагональная матрица с числами $\lambda_1, \dots, \lambda_N$ по диагонали. Такое представление матрицы называется спектральным разложением. Если положительная определенность Γ_k нарушена, то существует $\lambda_i \leq 0$. Матрица $\bar{\Gamma}_k$ строится так, что у нее сохраняются все собственные векторы u_1, \dots, u_N , а собственные числа заменяются на новые $\bar{\lambda}_i$ так, что

$$\bar{\lambda}_i = \begin{cases} \lambda_i, & \text{при } \lambda_i > \delta \\ \delta, & \text{при } \lambda_i \leq \delta, \end{cases} \quad (7.24)$$

где δ – малое положительное число. После этого полагается

$$\bar{\Gamma}_k = U \bar{L} U^T, \quad (7.25)$$

где в диагональной матрице \bar{L} на диагонали используются числа $\bar{\lambda}_1, \dots, \bar{\lambda}_N$.

При этом подпространство локальной положительной кривизны функции $f(y)$ сохраняется, а подпространство отрицательной кривизны функции заменяется подпространством малой положительной кривизны. Если построить квадратичную аппроксимацию функции $f(y)$ по результатам ее испытания в точке y^k и затем заменить в ней матрицу вторых производных Γ_k на модифицированную матрицу (7.24)–(7.25), то произойдут качественные изменения в структуре аппроксимации.

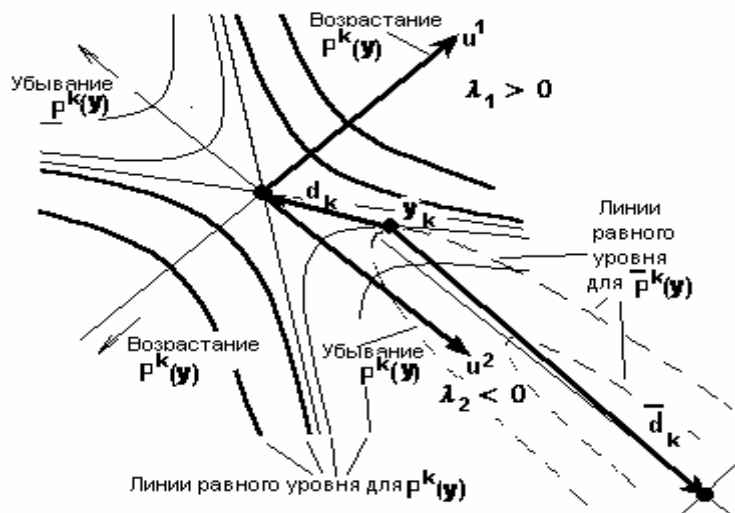


Рис. 7. 8. Изменение изолиний при замене матрицы

На рис.7.8 на примере пространства двух переменных показаны изменения в линиях равного уровня квадратичной аппроксимации $P^k(y)$ после замещения знаконеопределенной матрицы Γ_k положительно определенной $\bar{\Gamma}_k$. На этом рисунке видно, как направление метода Ньютона d^k , приводящее в стационарную точку поверхности $P^k(x)$ заменяется новым направлением \bar{d}^k , приводящим в минимум модифицированной аппроксимации. Для создания наглядных представлений об изменении характера аппроксимирующей поверхности при замене Γ_k на положительно определенную $\bar{\Gamma}_k$, полезно обратиться к иллюстрациям, представленным на рис. 7.8–7.9.

Поведение функции $f(x)$, представленное на рис.7.9, соответствует изолиниям, показанным на рис.7.8. На рисунке следует обратить внимание на то, что кривизна модифицированной аппроксимирующей поверхности $\bar{P}^k(y)$, построенной по измененной матрице $\bar{\Gamma}_k$, рассматриваемая в направлении u^1 положительной локальной кривизны поверхности $f(y)$, совпадает с локальной кривизной самой $f(y)$ в этом направлении. В то же время, в направлении u^2 отрицательная кривизна заменена малой положительной.

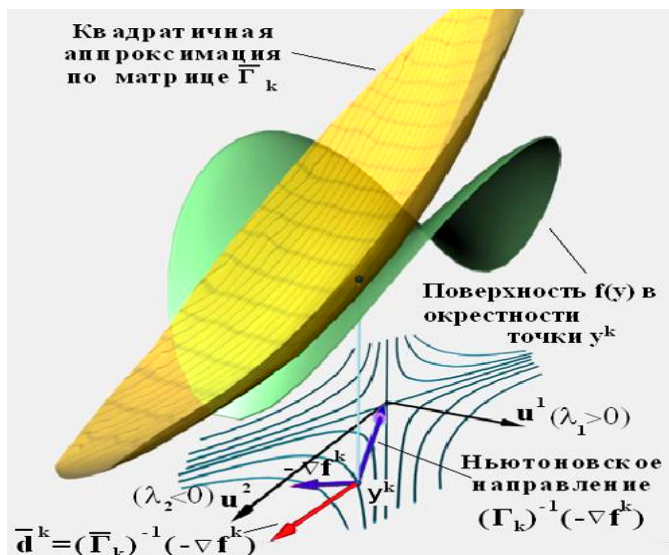


Рис. 7. 9. Изменение вида аппроксимирующей поверхности в случае знаконеопределенной матрицы Гессе

Несколько иное соответствие между первоначальной и модифицированной аппроксимациями возникает в том случае, когда матрица Γ_k отрицательно определена. В этом случае все собственные направления матрицы Γ_k трансформируются в направления малой положительной кривизны (рис.7.10). Это приводит к замене направления Ньютона d^k на новое — \bar{d}^k , ориентированное на точку минимума измененной аппроксимации.

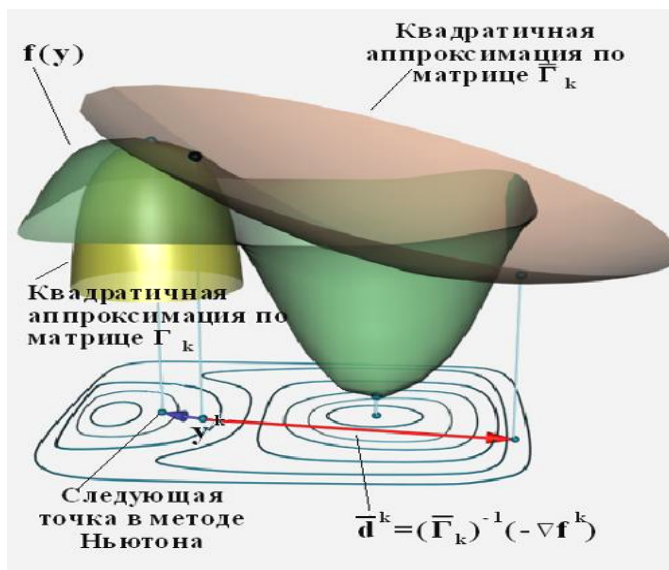


Рис. 7. 10. Влияние модификации матрицы при ее отрицательной определенности

Таким образом, метод коррекции матрицы на основе спектрального разложения весьма нагляден и прост для понимания. Однако этот подход имеет один существенный недостаток – большие затраты по вычислениям, связанные с поиском собственных векторов и чисел симметричной матрицы. Необходимый объем вычислений для построения спектрального разложения оценивается как $4N^3$.

При разработке вычислительных методов оптимизации для построения положительно определенных матриц $\bar{\Gamma}_k$ по исходным матрицам Γ_k вместо

спектрального разложения часто используется модифицированное разложение Холецкого [4], вычислительная реализация которого проще и требует меньшего числа операций.

В теории матриц известно, что для любой симметричной положительно определенной матрицы Γ существует нижняя треугольная матрица L с единичной диагональю и диагональная матрица D с положительной диагональю, что справедливо *разложение Холецкого* $\Gamma = LDL^T$, т.е.

$$\begin{pmatrix} \gamma_{11} & \gamma_{12} & \dots & \gamma_{1N} \\ \gamma_{12} & \gamma_{22} & \dots & \gamma_{2N} \\ \cdot & \cdot & \dots & \cdot \\ \gamma_{1N} & \gamma_{2N} & \dots & \gamma_{NN} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{12} & 1 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ l_{1N} & l_{2N} & \dots & 1 \end{pmatrix} \begin{pmatrix} d_{11} & 0 & \dots & 0 \\ 0 & d_{22} & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & d_{NN} \end{pmatrix} \begin{pmatrix} 1 & l_{12} & \dots & l_{1N} \\ 0 & 1 & \dots & l_{2N} \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 1 \end{pmatrix}, \quad (7.26)$$

где $d_{11}, \dots, d_{NN} > 0$.

Отсюда вытекает, что

$$\gamma_{11} = d_{11}, \quad \gamma_{1j} = d_{11}l_{1j}, \quad (j > 1) \quad (7.27)$$

$$\gamma_{ii} = d_{ii} + (l_{1i})^2 d_{11} + (l_{2i})^2 d_{22} + \dots + (l_{(i-1)i})^2 d_{(i-1)(i-1)}, \quad (i = 2, \dots, N) \quad (7.28)$$

$$\gamma_{ij} = (l_{1i})(l_{1j})d_{11} + (l_{2i})(l_{2j})d_{22} + \dots + (l_{(i-1)i})(l_{(i-1)j})d_{(i-1)(i-1)} + (l_{ij})d_{ii}, \quad (7.29)$$

для $i < j \leq N$.

Из (7.27)–(7.29) легко получить формулы для определения коэффициентов разложения Холецкого. Расчет их значений производится строка за строкой для матриц D и L^T . Порядок вычисления этих коэффициентов удобно пояснить следующей диаграммой:

$$d_{11} \Rightarrow l_{12}, l_{13}, l_{14}, \dots; \quad d_{22} \Rightarrow l_{23}, l_{24}, \dots; \quad d_{33} \Rightarrow l_{34}, l_{35}, \dots$$

Получим

$$d_{11} = \gamma_{11}, \quad l_{1j} = (1/d_{11})\gamma_{1j}, \quad (j = 2, \dots, N) \quad (7.30)$$

$$d_{ii} = \gamma_{ii} - \sum_{s=1}^{i-1} (l_{si})^2 d_{ss}, \quad (i > 1)$$

$$l_{ij} = \left(\gamma_{ij} - \sum_{s=1}^{i-1} (l_{si})(l_{sj})d_{ss} \right) / d_{ii}, \quad (i > 1, j = i+1, \dots, N) \quad (7.31)$$

Построим теперь *модифицированное разложение Холецкого* для произвольной (не обязательно положительно определенной) симметричной матрицы Γ . В процессе разложения будем производить коррекцию получаемых элементов d_{ij} так, чтобы модифицированные элементы \bar{d}_{ij} удовлетворяли условию

$$\bar{d}_{ij} \geq \delta > 0. \quad (7.32)$$

Это обеспечит положительность с «запасом» элементов модифицированной диагональной матрицы \bar{D} . Отметим, что условие (7.32) не может являться единственным условием модификации. Действительно, близость к нулю некоторых модифицированных элементов \bar{d}_{ij} при их дальнейшем использовании в (7.31) может привести к лавинообразному росту элементов l_{ij} при вычислениях. При обычном разложении Холецкого это невозможно, т.к. из (7.28) вытекает, что

$$(l_{si})^2 d_{ss} \leq \gamma_{ii} \leq \max\{\gamma_{ii} : i = 1, \dots, N\}, \quad (s = 1, \dots, i).$$

Обозначим через $\gamma^* = \max\{\gamma_{ii} : i = 1, \dots, N\}$ и введем $\beta^2 \geq \gamma^*$. Наложим требование, чтобы для модифицированных элементов разложения выполнялось

$$d_{ss} \leq \beta^2; ((\bar{l}_{si})^2 \bar{d}_{ss}) \leq \beta^2 (s = 1, \dots, N, i > s) \quad (7.33)$$

Выполним необходимую модификацию за счет изменения только диагональных элементов матрицы Γ . Обозначим через Δ_i добавки к элементам γ_{ii} . Тогда согласно (7.32)–(7.33) должно выполняться:

$$\bar{d}_{ii} = \gamma_{ii} + \Delta_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss} \geq \delta > 0,$$

$$\max_{j=i+1, \dots, N} (\bar{l}_{ij})^2 \bar{d}_{ii} = \left(\left(\max_{j=i+1, \dots, N} \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right)^2 \right) / \left(\gamma_{ii} + \Delta_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss} \right) \right) \leq \beta^2.$$

Если обозначить

$$c_i^2 = \max_{j=i+1, \dots, N} \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right)^2 \quad (7.34)$$

$$\tilde{d}_{ii} = \gamma_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss}, \quad (7.35)$$

то

$$\Delta_i = \max \{0; c_i^2 / \beta^2 - \tilde{d}_{ii}; \delta - \tilde{d}_{ii}\}.$$

Это соответствует выбору

$$\bar{d}_{ii} = \max \{ \tilde{d}_{ii}; c_i^2 / \beta^2; \delta \} \quad (7.36)$$

$$\bar{l}_{ij} = \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right) / \bar{d}_{ii} \quad (7.37)$$

Элементы (7.36), (7.37) определяют модифицированные матрицы \bar{L}^T , \bar{D} . В качестве положительно определенной матрицы – приближения для Γ используется $\bar{G} = \bar{L} \bar{D} \bar{L}^T$.

Заметим, что сумма квадратов поправок к элементам матрицы Γ равна $\Delta_1^2 + \dots + \Delta_N^2$. Для уменьшения этой величины в [4] рекомендуется выбирать

$$\beta^2 = \max \{ \gamma^*; \xi / (N^2 - 1)^{1/2}; \varepsilon_M \}, \quad (7.38)$$

где ε_M — наименьшее положительное вещественное число в машинной арифметике, а

$$\gamma^* = \max \{ \gamma_{ii} : i = 1, \dots, N \}, \quad \xi = \max \{ |\gamma_{ij}| : 1 \leq i \leq N, 1 \leq j \leq N, i \neq j \}. \quad (7.39)$$

Приведем пошаговое описание метода Ньютона с регулировкой шага и модификацией матрицы вторых производных на положительную определенность с использованием модифицированного преобразования Холесского.

АЛГОРИТМ МОДИФИЦИРОВАННОГО МЕТОДА НЬЮТОНА–РАФСОНА.

ШАГ 0. Задаются начальная точка y^0 , параметры выбора коэффициента одномерного шага $0 < \mu < \eta \ll 1$, $0 < \sigma \ll 1$; параметр останова ε и параметр модификации $\delta > 0$. Полагается $k = 0$.

ШАГ 1. Вычисляются $f^k = f(y^k)$, $\nabla f^k = \nabla f(y^k)$, $\Gamma_k = \Gamma^f(y^k)$.

ШАГ 2. Вычисляется β^2 по формулам (7.38), (7.39). Строятся матрицы модифицированного разложения Холесского \bar{L}_k , \bar{D}_k для матрицы Γ_k .


ШАГ 3. Определяется модифицированное направление Ньютоновского шага $\bar{d}^k = (\bar{\Gamma}_k)^{-1} (-\nabla f^k)$ путем последовательного решения двух систем линейных уравнений с треугольными матрицами


$$\begin{aligned} \bar{L}_k v^k &= -\nabla f^k \\ (\bar{D}_k (\bar{L}_k)^T) \bar{d}^k &= v^k \end{aligned}$$

ШАГ 4. Определяется x^k по алгоритму выбора коэффициента одномерного шага $x^k \in P$ из (7.13). Определяется $y^{k+1} = y^k + x^k d^k$

ШАГ 5. Вычисляется $f^{k+1} = f(y^{k+1})$, $\nabla f^{k+1} = \nabla f(y^{k+1})$, $\Gamma^{k+1} = \Gamma^f(y^{k+1})$, полагается $k := k+1$.

ШАГ 6. Если $\|\nabla f^k\| \leq \varepsilon$, то производится останов. Точка y^k выдается в качестве оценки решения. Если же $\|\nabla f^k\| \geq \varepsilon$, то осуществляется переход на шаг 2.

 **Замечание 1.** Известно, что выполнение модифицированного разложения Холесского требует около $(1/6)N^3$ операций, а последующее определение d^k требует числа операций порядка N^2 .

 **Замечание 2.** Модифицированный метод Ньютона сохраняет поисковые возможности на функциях с областями плохого поведения (вырожденность, знаконеопределенность матриц Γ_k), поскольку новое направление поиска \bar{d}^k , в силу гарантированной положительной определенности матриц $\bar{\Gamma}_k$, обязательно является направлением строгого локального убывания и соответствует направлению антиградиента в некоторой новой метрике пространства. Кроме того, метод обладает сверхлинейной скоростью сходимости в окрестности решения, если функция в этой окрестности дважды непрерывно дифференцируема и принадлежит классу $\Phi_{m,M}$ из (7.22).

Справедливость последнего утверждения следует из того, что в указанной области модификация матрицы производиться не будет, т.е. $\bar{\Gamma}_k = \Gamma_k$, а, следовательно, метод будет точно совпадать с методом Ньютона–Рафсона.

7.5. Методы первого порядка, явно изменяющие метрику пространства

В этом разделе будет продолжено изучение методов, основанных на квадратичной модели поведения минимизируемой функции. В отличие от предыдущего раздела, будет рассмотрена группа методов, которые хотя и используют предположение о гладкости функции и близости ее к квадратичной, но не измеряют матриц вторых производных. Будет рассмотрено несколько групп таких методов. Методы первой группы (квазиньютоновские методы) строят оценки матриц Гессе и используют их вместо истинных матриц вторых производных, точно также, как это делает метод Ньютона–Рафсона. Методы второй группы — методы растяжения также основаны на построении вспомогательных матриц, используемых для перемасштабирования пространства, но эти матрицы не являются оценками Гессеана и строятся на основе эвристических принципов. Эти методы будут рассмотрены в данном разделе.

Методы третьей группы явно никаких матриц не строят, хотя неявно метрику пространства изменяют (методы сопряженных направлений). Они рассматриваются в разделе 7.6.

7.5.1. Квазиньютоновские методы. Рекуррентные соотношения для оценок матриц Гессе по измерениям градиента в основных точках поиска

Методы этого класса относятся к методам первого порядка. Они используют результаты испытаний, состоящих в вычислении $f(y^k)$, $\nabla f(y^k)$. Предполагается, что функция f обладает свойствами, соответствующими квадратичной модели.

Следовательно, у функции f существует симметричная матрица вторых производных, недоступная непосредственному измерению.

Казалось бы, в этих условиях самым естественным являлось конечно-разностное оценивание Гессииана в каждой точке y^k по измерениям градиента на множестве узлов, размещенных с некоторым шагом h в окрестности данной точки. Получив оценку можно применить модифицированный метод Ньютона–Рафсона. Однако данный подход требует слишком большого объема вычислений и, кроме того, связан со значительными погрешностями оценивания. Это приводит к тому, что данный подход обычно не применяется. Оказывается оценки Гессииана можно строить без дополнительных вычислений градиента функции $f(y)$. Именно такой подход используется в квазиньютоновских методах.

Идея, положенная в основу *квазиньютоновских методов*, состоит в том, чтобы по результатам измерения градиентов функции f в точках y^k траектории поиска попытаться построить симметричную матрицу \mathbf{G}_k , являющуюся оценкой кривизны поверхности $f(y)$ на траектории поиска, т.е. оценивающую $\mathbf{\Gamma}^f$, или построить симметричную матрицу \mathbf{H}_k , являющуюся оценкой обратной матрицы $(\mathbf{\Gamma}^f)^{-1}$. После выполнения k -го испытания, \mathbf{G}_k рассматривается как оценка матрицы вторых производных $\mathbf{\Gamma}^f(y^k)$, а \mathbf{H}_k — как оценка $(\mathbf{\Gamma}^f)^{-1}$. В зависимости от того, какая именно матричная оценка используется, точка очередного измерения будет выбираться по правилу

$$y^{k+1} = y^k + x^k d^k, \quad (7.40)$$

$$d^k = (\mathbf{G}_k)^{-1}(-\nabla f^k) \quad (\text{или } d^k = (\mathbf{H}_k)(-\nabla f^k)), \quad (7.41)$$

$$x^k \in \Pi = \Pi_1(\eta) \cap \Pi_2(\mu), \quad 0 < \mu < \eta < 1.$$

Методы вида (7.40), (7.41) выбирают направления перемещения, совпадающие с антиградиентными направлениями функции f , вычисленными в некотором измененном пространстве. Фактически, на каждом шаге выполняется изменение метрики пространства переменных (поворот осей и перемасштабирование) за счет домножения градиента на соответствующую матрицу. Поэтому второе название этих методов — *методы переменной метрики*. При определенных условиях через N шагов матрица \mathbf{G}_N будет близка к матрице $\mathbf{\Gamma}^f(y^N)$, а \mathbf{H}_N — к $(\mathbf{\Gamma}^f(y^N))^{-1}$.

Рассмотрим основные идеи, лежащие в основе построения оценок \mathbf{G}_k и \mathbf{H}_k . Рассмотрим случай, когда $f(y) = (y^T \mathbf{\Gamma} y)/2 + (c, y) + b$ — квадратичная функция с симметричной положительно определенной матрицей $\mathbf{\Gamma}$.

Пусть

$$\begin{aligned} \Delta^k &= y^{k+1} - y^k, \\ z^k &= \nabla f^{k+1} - \nabla f^k, \end{aligned} \quad (7.42)$$

где $\nabla f^k = \nabla f(y^k)$, $\nabla f^{k+1} = \nabla f(y^{k+1})$. Тогда, очевидно, будут выполняться равенства

$$z^k = \nabla f(y^k + \Delta^k) - \nabla f(y^k) = \mathbf{\Gamma} \Delta^k, \quad (\mathbf{\Gamma}^{-1}) z^k = \Delta^k.$$

Потребуем, чтобы таким же условиям удовлетворяли, соответственно, оценка \mathbf{G}_{k+1} матрицы $\mathbf{\Gamma}$ и оценка \mathbf{H}_k матрицы $\mathbf{\Gamma}^{-1}$, построенные по $(k+1)$ -му измерению градиента. А именно, пусть выполняется требование

$$z^k = \mathbf{G}_{k+1} \Delta^k \quad (\mathbf{H}_{k+1} z^k = \Delta^k). \quad (7.43)$$

Условия (7.43) называются *квазиньютоновскими условиями*.

Наложим дополнительные требования на матрицы оценок. Поскольку сама матрица \mathbf{G} симметрична, потребуем выполнения свойства симметрии от матрицы \mathbf{G}_{k+1} , (\mathbf{H}_{k+1}), положив

$$\mathbf{G}_{k+1} = (\mathbf{G}_{k+1})^T \quad (\text{или } \mathbf{H}_{k+1} = (\mathbf{H}_{k+1})^T). \quad (7.44)$$

Будем определять ее новое значение путем коррекции предыдущей матрицы

$$\mathbf{G}_{k+1} = \mathbf{G}_k + \mathbf{U}_k \quad (\text{или } \mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{U}_k).$$

где поправки \mathbf{U}^k строятся в виде матриц ранга 1 и находятся из условий (7.43), (7.44).

Эти условия определяют поправку неединственным образом. Рассмотрим возможные приемы ее построения на примере оценок \mathbf{G}_k . Оценки вида \mathbf{H}_k строятся аналогично.

Простейший способ определения поправочной матрицы, предложенный Бройденом, состоит в том, чтобы составить ее из вектор–столбцов невязок вида $z^k - \mathbf{G}_k \Delta^k$, помноженных на специально подобранные числа. Нетрудно проверить, что при произвольных v^k , для которых $(v^k)^T \Delta^k \neq 0$, нужная оценка определяется следующей формулой

$$\mathbf{G}_{k+1} = \mathbf{G}_k + (z^k - \mathbf{G}_k \Delta^k)(v^k)^T / ((v^k)^T \Delta^k). \quad (7.45)$$

Если положить в ней $v^k = (z^k - \mathbf{G}_k \Delta^k)$, то получим *формулу Бройдена (B–формулу)*

$$\mathbf{G}_{k+1} = \mathbf{G}_k + (z^k - \mathbf{G}_k \Delta^k)(z^k - \mathbf{G}_k \Delta^k)^T / ((z^k - \mathbf{G}_k \Delta^k)^T \Delta^k) \quad (7.46)$$

Существуют и другие способы оценивания. Например, можно несимметричную поправку в формуле (7.45) заменить похожей симметричной. Непосредственной проверкой можно убедиться (полезно в качестве упражнения выполнить эту проверку), что для любого вектора v^k такого, что $(v^k)^T \Delta^k \neq 0$, соотношение

$$\begin{aligned} \mathbf{G}_{k+1} = & \mathbf{G}_k + ((z^k - \mathbf{G}_k \Delta^k)(v^k)^T + v^k(z^k - \mathbf{G}_k \Delta^k)^T) / ((v^k)^T \Delta^k) - \\ & - (v^k)(v^k)^T (z^k - \mathbf{G}_k \Delta^k)^T \Delta^k / ((v^k)^T \Delta^k)^2 \end{aligned} \quad (7.47)$$

дает оценочную матрицу, удовлетворяющую (7.43), (7.44).

Положив в (7.45) $v^k = z^k$, получим *формулу Девидона–Флетчера–Пауэлла (DFP–формулу)*, представимую в следующем виде

$$\mathbf{G}_{k+1} = \mathbf{G}_k - \mathbf{G}_k \Delta^k (\Delta^k)^T \mathbf{G}_k / ((\Delta^k)^T \mathbf{G}_k \Delta^k) + z^k (z^k)^T / ((z^k)^T \Delta^k) + (\Delta^k)^T \mathbf{G}_k \Delta^k w^k (w^k)^T, \quad (7.48)$$

где

$$w^k = z^k / ((z^k)^T \Delta^k) - \mathbf{G}_k \Delta^k / ((\Delta^k)^T \mathbf{G}_k \Delta^k). \quad (7.49)$$

Поскольку в (7.48), (7.49) $w^k (w^k)^T$ – симметричная матрица и, как можно показать, $(w^k)^T \Delta^k = 0$ (проверьте это в качестве упражнения), то, в силу (7.43), последнее слагаемое в (7.48) можно отбросить. В результате будет получена *сокращенная формула Девидона–Флетчера–Пауэлла*. Ее называют также *формулой Бройдена–Флетчера–Гольдфарба–Шанно (BFGH–формула)*.

В приведенных итерационных соотношениях начальное значение выбирается в виде $\mathbf{G}_0 = \mathbf{E}$ (единичная матрица).

Аналогичные формулы для оценивания обратной матрицы Гессе имеют следующий вид.

B – формула для \mathbf{H}_k

$$\mathbf{H}_{k+1} = \mathbf{H}_k + (\Delta^k - \mathbf{H}_k z^k)(\Delta^k - \mathbf{H}_k z^k)^T / ((\Delta^k - \mathbf{H}_k z^k)^T z^k) \quad (7.46')$$

Сокращенная DFP – формула для \mathbf{H}_k

$$\mathbf{H}_{k+1} = \mathbf{H}_k - \mathbf{H}_k \mathbf{z}^k (\mathbf{z}^k)^T \mathbf{H}_k / ((\mathbf{z}^k)^T \mathbf{H}_k \mathbf{z}^k) + \Delta^k (\Delta^k)^T / ((\Delta^k)^T \mathbf{z}^k), \quad (7.48')$$


В теории оптимизации известно удивительное свойство описанных выше матричных оценок. Сформулируем его в виде теоремы.

Теорема 7.5. Для квадратичных функций $f(x)$, $x \in \mathbb{R}^N$ с положительно определенными матрицами вторых производных матрица G_N (H_N), полученная с использованием процедур (7.40)–(7.42), а также B, DFP или сокращенной DFP– формул (7.46), (7.48) ((7.46'), (7.48')), будет совпадать с матрицей вторых производных Γ^f (с обратной матрицей $(\Gamma^f)^{-1}$) функции f , а точка y^N — с глобальным минимумом y^* функции f . Кроме того, $\forall k \leq N$ матрицы G_k (H_k) будут симметричны и положительно определены.

Доказательство этого факта можно найти в [1] применительно к сокращенной DBF–формуле для \mathbf{H}_k . Оно не является сложным, но достаточно громоздко, включает несколько этапов, которые обосновываются по индукции. Здесь это доказательство опускается.

Построенные алгоритмы могут быть применены для достаточно произвольных функций f , не являющихся квадратичными. В этом случае обычно $\mathbf{G}_N \neq \Gamma^f(y_N)$, поскольку матрица вторых производных не постоянна (то же относится и к оценке обратной матрицы). После каждой серии из N шагов необходим повторный запуск метода из получаемой точки y^N . Кроме того, процесс поиска может привести к тому, что матрицы \mathbf{G}_k могут оказаться вырожденными или знаконеопределенными. При этом направление шага d^k в (7.41) перестанет быть направлением убывания функции и величина смещения x^k в (7.40) окажется равной нулю. Простейший способ коррекции в этом случае состоит в замене направления, построенного по правилу (7.41), на обычное антиградиентное направление (хотя эта стратегия не является лучшей).

Следует отметить, что в неквадратичном случае безопаснее использовать оценки прямой матрицы, а не обратной. Стратегия поиска в квазиньютоновских методах проиллюстрирована на рис.7.11. На этом рисунке показан (пунктиром) возможный вид линий равного уровня для квадратичных аппроксимаций функции $f(y)$, построенных по матричным оценкам \mathbf{G}_k . Поскольку $\mathbf{G}_0 = \mathbf{E}$, то первая из этих линий уровня, построенная для точки y^0 , является окружностью, а на остальных шагах окружности преобразуются в эллипсы. Выбираемые далее методом направления поиска d^k проходят через центры этих эллипсов, являющиеся стационарными точками построенных квадратичных аппроксимаций. Эти направления являются антиградиентными направлениями в пространствах с новыми метриками, связанными с матрицами \mathbf{G}_k . Они могут сильно отличаться от направлений градиента в исходном пространстве.

 **Замечание.** При реализации алгоритма требуется проверка положительной определенности матриц \mathbf{G}_k , а также вычисление направления поиска согласно (7.41), т.е. определение $d^k = (\mathbf{G}_k)^{-1}(-\nabla f^k)$. Эти операции можно выполнить совместно. Для этого достаточно выполнить для матрицы \mathbf{G}_k разложение Холецкого. Если при этом для диагональных элементов матрицы \mathbf{D}_k нарушится условие положительности $d_{ii} > 0$, то в качестве направления поиска нужно выбрать обычное антиградиентное направление, если же разложение $\mathbf{G}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^T$ будет построено, то определение вектора d^k сведется к решению двух линейных систем с треугольными матрицами $\mathbf{L}_k v = -\nabla f^k$ и $\mathbf{D}_k (\mathbf{L}_k)^T d^k = v$.

Линии равного уровня функции $f(y)$

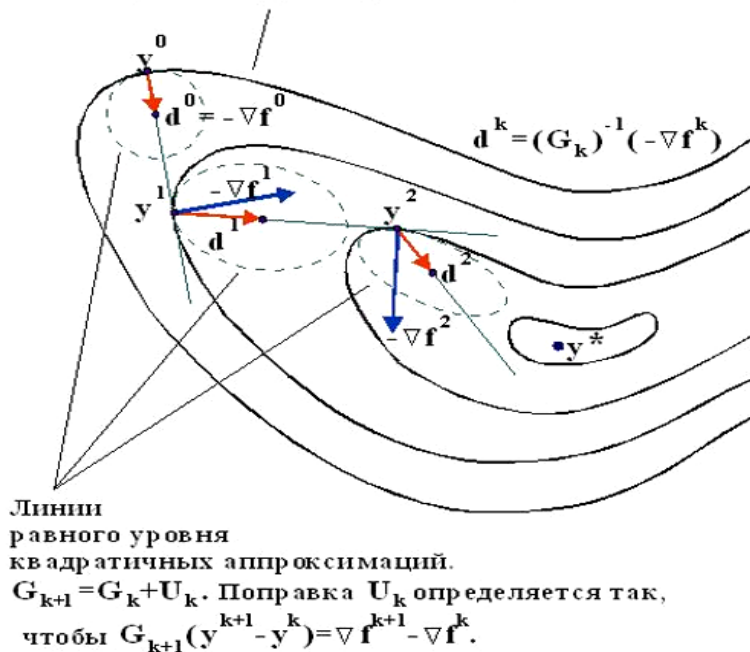


Рис. 7.11. Направления поиска в квазиньютоновских методах

ОПИСАНИЕ АЛГОРИТМА

ШАГ 0. Определяем $\varepsilon > 0$ — параметр останова, μ , η и σ — параметры одномерного поиска ($0 < \mu < \eta < 1$, $0 < \sigma < 1$). Задаем точку начала поиска y^0 .

ШАГ 1. Полагаем $G_0 = E$ и вычисляем $f^0 = f(y^0)$, $\nabla f^0 = \nabla f(y^0)$, $k = 0$.

ШАГ 2. Выполняем преобразование Холесского для матрицы G_k . Если преобразование выполнить не удалось, полагаем $d^k = (-\nabla f^k)$ и переходим на шаг 4. В противном случае получаем $G_k = L_k D_k L_k^T$.

ШАГ 3. Определяем направление поиска $d^k = (G_k)^{-1}(-\nabla f^k)$ путем решения двух систем с треугольными матрицами

$$\begin{aligned} L_k v &= -\nabla f^k \\ D_k (L_k)^T d^k &= v \end{aligned} \quad (7.50)$$

ШАГ 4. Определяем $x^k \in \Pi$ с помощью алгоритма выбора одномерного шага, вычисляем

$$\begin{aligned} y^{k+1} &= y^k + x^k d^k \\ f^{k+1} &= f(y^{k+1}), \nabla f^{k+1} = \nabla f(y^{k+1}) \\ \Delta^k &= y^{k+1} - y^k, z^k = \nabla f^{k+1} - \nabla f^k. \end{aligned}$$

ШАГ 5. Если $k=N$, проверяем критерий останова: при $\|\nabla f^{k+1}\| \leq \varepsilon$ останавливаем поиск и принимаем y^{k+1} в качестве решения; при $\|\nabla f^{k+1}\| > \varepsilon$ полагаем $y^0 = y^{k+1}$ и переходим к шагу 1. Если $k \neq N$, то полагаем $k = k+1$ и переходим к шагу 2.

ШАГ 6. Производим вычисление матрицы G_{k+1} по B -формуле (7.46), DFP -формуле (7.48) или по $BFGH$ -формуле. Переходим на шаг 2.

7.5.2. Модифицированные квазиньютоновские методы

Материал излагается на примере методов, использующих оценки G_k прямой матрицы Гессе. В этих методах в случаях нарушения положительной определенности оценочных матриц G_k вместо использования антиградиентного направления выполняется замена матрицы G_k на близкую к ней положительно определенную матрицу \bar{G}_k , построенную с использованием модифицированного преобразования Холесского (или другого подобного преобразования). Рис.7.12 показывает изменения вида квадратичной аппроксимации функции $f(y)$ в результате выполнения указанного преобразования для случаев различной знакоопределенности матрицы G_k .

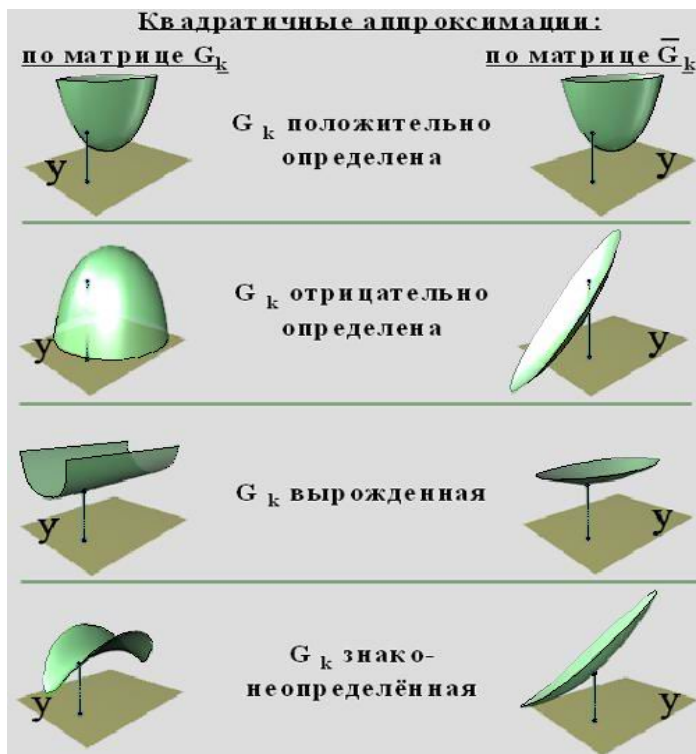


Рис. 7.12. Влияние модификации матрицы G_k на вид аппроксимирующей поверхности

ОПИСАНИЕ АЛГОРИТМА.

ШАГ 0. Определяем $\varepsilon > 0$ — параметр останова, δ — параметр модификации матрицы в модифицированном преобразовании Холесского ($\delta > 0$), μ , η и σ — параметры одномерного поиска ($0 < \mu < \eta < 1$, $0 < \sigma < 1$). Задаем точку начала поиска y^0 .

ШАГ 1. Полагаем $G_0 = E$ и вычисляем $f^0 = f(y^0)$, $\nabla f^0 = \nabla f(y^0)$, $k = 0$.

ШАГ 2. Выполняем модифицированное преобразование Холесского для матрицы G_k , получаем $G_k \Rightarrow \bar{G}_k = \bar{L}_k \bar{D}_k \bar{L}_k^T$.

ШАГ 3. Определяем направление поиска $d^k = (\bar{G}_k)^{-1}(-\nabla f^k)$ путем решения двух систем с треугольными матрицами

$$\begin{aligned} \bar{L}_k v &= -\nabla f^k \\ \bar{D}_k (\bar{L}_k)^T d^k &= v \end{aligned}$$

ШАГ 4. Определяем $x^k \in P$ с помощью алгоритма одномерного шага, вычисляем

$$\begin{aligned} y^{k+1} &= y^k + x^k d^k, \\ f^{k+1} &= f(y^{k+1}), \nabla f^{k+1} = \nabla f(y^{k+1}), \\ \Delta_k &= y^{k+1} - y^k, z^k = \nabla f^{k+1} - \nabla f^k. \end{aligned}$$

ШАГ 5. Если $k=N$, проверяем критерий останова: при $\|\nabla f^{k+1}\| \leq \varepsilon$ останавливаем поиск и принимаем y^{k+1} в качестве решения; при $\|\nabla f^{k+1}\| > \varepsilon$ полагаем $y^0 = y^{k+1}$ и переходим к шагу 1. Если $k \neq N$, то полагаем $k = k+1$ и переходим к шагу 6.

ШАГ 6. Производим вычисление матрицы G_{k+1} по B -формуле (7.46), DFP -формуле (7.48) или по $BFGH$ -формуле. Переходим на шаг 2.

Практический опыт показывает, что для широкого класса гладких задач описанные алгоритмы достаточно экономичны по числу шагов.

7.5.3. Методы растяжения пространства (R-алгоритмы Н.З. Шора)

Автором этой группы методов является украинский математик Н.З. Шор. Предложенные им методы основаны на подборе матрицы преобразования пространства. Для класса выпуклых гладких функций методы растяжения тесно связаны с методами А.С. Немировского и Д.Б. Юдина, рассмотренными в разделе 7.1 Преобразования пространства в алгоритмах Шора сводятся к последовательным растяжениям в специально подбираемых направлениях. Эти методы называют *R-алгоритмами*. Они по структуре близки к квазиньютоновским методам переменной метрики, но основаны не на оценке матрицы вторых производных, а на построении матрицы преобразования B_k , определяющей возврат от некоторых новых координат z к исходным: $y = B_k z$. Матрица B_k строится как произведение матриц преобразования $R_{\beta}(\xi^e)$, выполняющих растяжение или сжатие пространства z в β раз в направлениях ξ^e ($e = 1, 2, \dots, k$), $\|\xi^e\| = 1$.

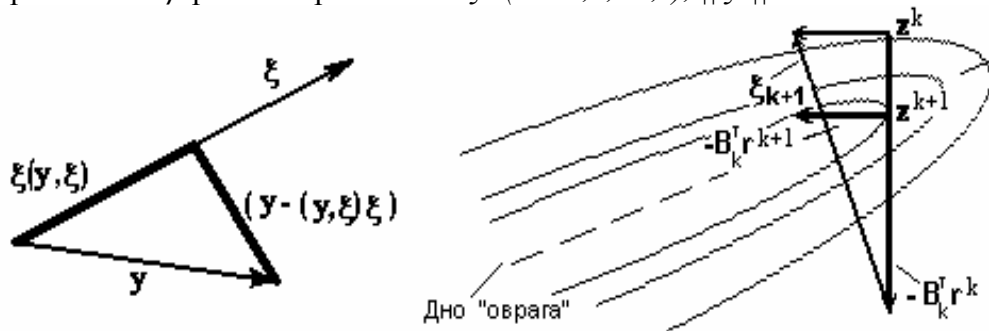


Рис. 7.13. Стратегия построения матрицы растяжения пространства

Нетрудно увидеть (рис. 7.13), что

$$R_{\beta}(\xi) y = (y - (y, \xi) \xi) + \beta (y, \xi) \xi = (E + (\beta - 1) \xi \xi^T) y. \quad (7.51)$$

Следовательно, $R_{\beta}(\xi) = E + (\beta - 1) \xi \xi^T$.

Пусть $r^k = \nabla f(y^k)$ – градиент функции в исходном пространстве, а \bar{r}^k — это значение градиента, подсчитанного в соответствующей точке z в новом пространстве переменных. Тогда

$$\bar{r}^k = \nabla_z f(B_k z^k) = B_k^T r^k.$$

Для отыскания минимума функции $f(y)$ будем использовать схему метода наискорейшего градиентного поиска, но так, чтобы на каждом шаге k градиент вычислялся в новом пространстве, связанном с матрицей преобразования B_k (рис.7.14).

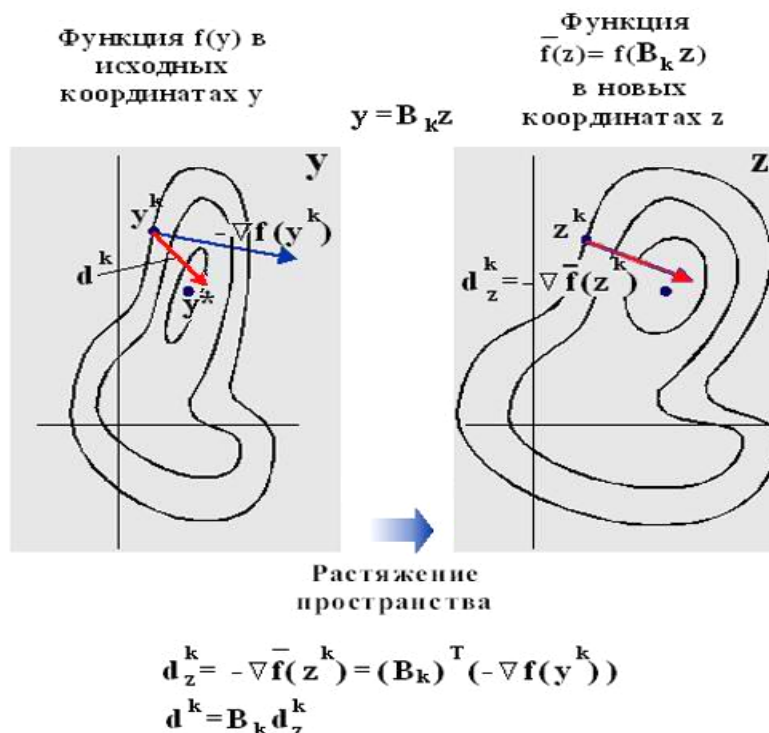


Рис. 7.14. Выбор направления в методе растяжения пространства

В этом пространстве будем в качестве очередного направления растяжения выбирать вектор $\xi^{k+1} = B_k^T(r^k - r^{k-1})$, определяющий разность двух последовательных измерений вектора градиента в пространстве, связанном с B_k . Этот вектор будет близок к нормали для многообразия, на котором лежит дно оврага минимизируемой функции (см. рис.7.13), если рассматривать эти объекты в пространстве новых переменных z .

В найденном направлении ξ^{k+1} будем осуществлять дополнительное растяжение пространства в фиксированное число раз (с коэффициентом $\alpha \approx 2$ или 3). При возврате к исходным координатам этой операции будет соответствовать сжатие в направлении ξ^{k+1} с коэффициентом $\beta = 1/\alpha$. Следовательно, $B_{k+1} = B_k R_{1/\alpha}(\xi^{k+1})$.

Мы приходим к следующему АЛГОРИТМУ МЕТОДА РАСТЯЖЕНИЯ.

ШАГ 0. Задаются $\varepsilon > 0$ — параметр критерия останова, $0 < \mu < \eta \ll 1$, $0 < \sigma \ll 1$ — параметры алгоритма выбора коэффициента одномерного шага, y^0 — начальная точка поиска, α — коэффициент растяжения пространства.

ШАГ 1. Вычисляются $f^0 = f(y^0)$, $r^0 = \nabla f(y^0)$, полагается $B_0 = E$, $k = 0$.

ШАГ 2. Вычисляется величина коэффициента одномерного шага x^k методом "аккуратного" одномерного поиска. Определяются

$$\begin{aligned}
 y^{k+1} &= y^k + x^k B_k (B_k)^T (-r^k), \\
 f^{k+1} &= f(y^{k+1}), \quad r^{k+1} = \nabla f(y^{k+1}).
 \end{aligned}
 \tag{7.52}$$

ШАГ 3. Если $\|r^{k+1}\| < \varepsilon$, то выполняется останов метода поиска, иначе переходим к шагу 4.

ШАГ 4. Выбирается направление дополнительного растяжения $\xi^{k+1} = (B_k)^T (r^{k+1} - r^k)$ и выполняется его нормировка $\xi^{k+1} := \xi^{k+1} / \|\xi^{k+1}\|$.

ШАГ 5. Пересчитывается матрица преобразования с учетом растяжения пространства в α раз вдоль ξ^{k+1} :

$$B_{k+1} = B_k R_{1/\alpha}(\xi^{k+1}). \quad (7.53)$$

ШАГ 6. Если хотя бы один из элементов b_{ij} матрицы B_{k+1} превысит по модулю некоторое заранее установленное пороговое значение, то все элементы этой матрицы делятся на модуль элемента b_{ij} . Изменяется $k=k+1$ и выполняется переход к шагу 2.

7.6. Методы сопряженных направлений

7.6.1. Сопряженные направления и их свойства

Построение методов *сопряженных направлений* основано на квадратичной модели поведения минимизируемой функции. Предположим, что $f(y)$ — квадратичная функция (7.15) с положительно определенной матрицей.

Определение 7.8. Система линейно-независимых векторов p^0, p^1, \dots, p^{N-1} для симметричной матрицы Γ называется Γ -сопряженной, если

$$\forall i=1, \dots, N; j=1, \dots, N; i \neq j: (p^i, \Gamma p^j) = 0. \quad (7.54)$$


Определение 7.9. Пусть M — линейное многообразие, Γ — симметричная матрица, $x \neq 0$ и $x \notin M$ и

$$\forall z \in M: (x, \Gamma z) = 0, \quad (7.55)$$

тогда вектор x называется Γ -сопряженным с многообразием M .

Можно легко доказать следующую лемму.

Лемма 7.2 Если p^0, p^1, \dots, p^{N-1} — все отличны от нуля, Γ — не только симметрична, но еще и положительно определена, тогда из (7.54) следует линейная независимость векторов p^0, p^1, \dots, p^{N-1} .

 **Замечание.** В условиях леммы сопряженность означает ортогональность в смысле некоторого нового скалярного произведения.

Для построения методов, использующих сопряженные направления, чрезвычайно важным является свойство, определяемое следующей леммой.

Лемма 7.3 Пусть $f(y)$ — квадратичная функция вида (7.16) с симметричной положительно определенной матрицей Γ , а M и L — линейные многообразия, причем $M \subset L$, тогда, если z — точка минимума $f(y)$ на M , а u — точка минимума $f(y)$ на L , то вектор $(u-z)$ будет Γ -сопряжен с многообразием M .

Это утверждение иллюстрируется на рис. 7.15.

Доказательство проведем следующим образом. Рассмотрим произвольный вектор $e \in M$. Поскольку $\nabla f(y) = \Gamma y + c$, а матрица Γ симметрична, то $(u-z, \Gamma e) = (\Gamma u - \Gamma z, e) = (\nabla f(u) - \nabla f(z), e)$. Последнее скалярное произведение равно нулю, т.к. по теореме Лагранжа в точках минимума u и z на линейных многообразиях L и M градиенты функции ортогональны этим многообразиям,

а поскольку $e \in M \subset L$, то $\forall e \in M: (\nabla f(u), e) = 0, (\nabla f(z), e) = 0$. Таким образом, для $x = u - z$ выполнено (7.55), следовательно $(u - z)$ будет Γ -сопряжен с многообразием M .

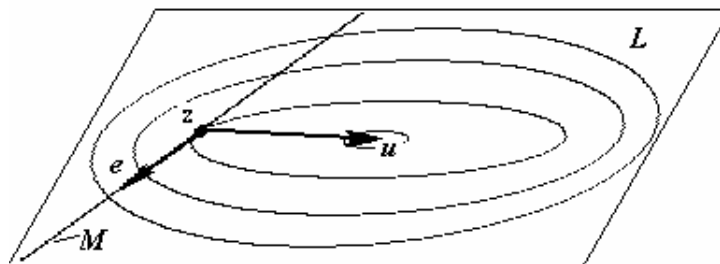


Рис. 7.15. Иллюстрация к лемме 7.3

Построим теперь вычислительные процедуры поиска минимума квадратичной функции $f(y)$, использующие Γ -сопряженные направления.

Определение 7.10. Поисковые процедуры вида (7.56)–(7.57) называются методами сопряженных направлений.

$$y^{k+1} = y^k + x^k p^k \quad (7.56)$$

$$f(y^k + x^k p^k) = \min\{f(y^k + x p^k): -\infty < x < +\infty\}. \quad (7.57)$$

Применение сопряженных направлений при построении методов оптимизации связано с замечательным свойством этих направлений приводить в минимум строго выпуклой квадратичной функции не более чем за N шагов.

Теорема 7.6. Пусть $f(y)$ — квадратичная функция вида (7.15) с симметричной положительно определенной матрицей Γ , а p^0, p^1, \dots, p^{N-1} — система Γ -сопряженных векторов. Тогда для любой начальной точки y^0 процедура поиска вида (7.56)–(7.57) приводит в минимум квадратичной функции с симметричной положительно определенной матрицей Γ ровно за N шагов, т.е. $y^N = y^*$, $f(y^N) = f(y^*)$.

Доказательство [10]. При поиске вдоль направления p^0 метод определит точку y^1 — минимум на одномерном многообразии $L(p^0)$, натянутом на p^0 . На втором шаге при поиске вдоль направления p^1 метод определит точку y^2 . По построению вектор $y^2 - y^1$ будет сопряжен с $L(p^0)$, т.е. ортогонален к p^0 в смысле нового скалярного произведения. Если теперь рассмотреть линейное многообразие $L(p^0, p^1)$, натянутое на p^0, p^1 и предположить, что минимум функции $f(y)$ достигается на нем в некоторой точке $\bar{y}^2 \neq y^2$, то возникнет противоречие. Действительно, по лемме 7.2 мы получим еще один вектор $\bar{y}^2 - y^1$, не принадлежащий прямой, проходящей через y^2 и y^1 , лежащий в том же двумерном многообразии $L(p^0, p^1)$ и ортогональный (в смысле нового скалярного произведения) к p^0 . Значит $\bar{y}^2 = y^2$, и на втором шаге метод сопряженных направлений найдет минимум на двумерном многообразии $L(p^0, p^1)$.

Продолжая аналогичные рассуждения приходим к выводу, что за N шагов метод найдет минимум на линейном многообразии $L(p^0, \dots, p^{N-1})$ размерности N , т.е. во всем пространстве (рис.7.16).

Для того, чтобы можно было воспользоваться методом сопряженных направлений необходим алгоритм вычисления Γ -сопряженных векторов p^0, \dots, p^{N-1} . Проблема, которая на первый взгляд кажется непреодолимой, заключается в том, чтобы построить Γ -сопряженные векторы не зная самой матрицы Γ . Однако,

как будет показано в следующем разделе, эта задача может быть решена с использованием результатов испытаний функции $f(y)$.

$$f(y) = y^T \Gamma y / 2 + c^T x$$

p^0, \dots, p^n - Γ -сопряжённые

$$(p^i, \Gamma p^j) = 0, \quad i \neq j$$

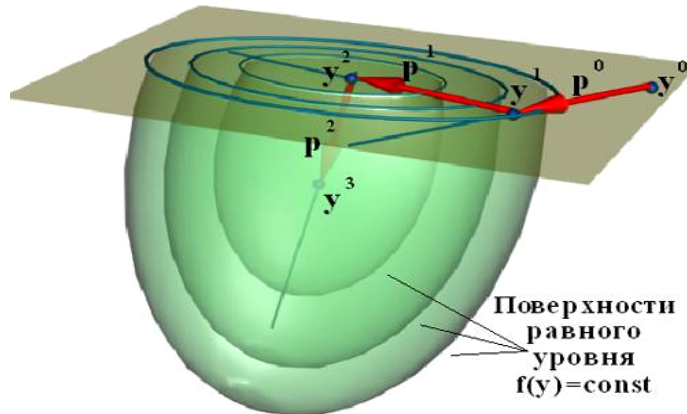


Рис. 7.16. Замечательное свойство сопряженных направлений

7.6.2. Метод сопряженных градиентов Флетчера-Ривса

Рассмотрим класс методов сопряженных направлений первого порядка, когда в результате испытания функции f в точке y^k определяются значения $f(y^k)$ и $\nabla f(y^k)$. Для построения метода сопряженных направлений необходимо по результатам испытаний построить систему Γ -сопряженных векторов p^0, \dots, p^{N-1} при условии, что сама матрица Γ является неизвестной.

Построим один из возможных методов такого типа – метод сопряженных градиентов Флетчера-Ривса (1964 год) [10]. Выберем

$$p^0 = -\nabla f^0, \quad \nabla f^0 = \nabla f(y^0). \quad (7.58)$$

Пусть векторы p^0, \dots, p^{k-1} построены. Положим

$$p^k = -\nabla f^k + \beta_{k-1} p^{k-1}, \quad \nabla f^k = \nabla f(y^k), \quad (7.59)$$

где y^k определяется условиями (7.56), (7.57). Подберем β_{k-1} из условия

$$(p^k, \Gamma p^{k-1}) = 0.$$

Получим

$$\beta_{k-1} = ((\nabla f^k)^T \Gamma p^{k-1}) / ((p^{k-1})^T \Gamma p^{k-1}) \quad (7.60)$$

Значение x^k , удовлетворяющее (7.57) для функции f , можно получить из условия экстремума $(\nabla f(y^k + x^k p^k), p^k) = 0$, если его переписать в виде $((\nabla f^k)^T p^k) + ((p^k)^T \Gamma p^k) x^k = 0$. Отсюда можно показать, что x^k будет иметь вид

$$x^k = -((p^k)^T \nabla f^k) / ((p^k)^T \Gamma p^k) = -((\nabla f^k)^T \nabla f^k) / ((\nabla f^k)^T \Gamma p^k). \quad (7.61)$$

Для этого в числителе и знаменателе первой дроби необходимо выразить p^k из (7.59) и воспользоваться тем, что $((p^{k-1})^T \nabla f^k) = 0$ по теореме Лагранжа, и $(p^{k-1}, \Gamma p^k) = 0$ по построению.

Кроме того, умножая (7.56) на Γ , получим дополнительное соотношение

$$\nabla f^{k+1} = \nabla f^k + x^k \Gamma p^k. \quad (7.62)$$

Лемма 7.4. Последовательность векторов градиентов $\nabla f^0, \nabla f^1, \dots, \nabla f^{N-1}$ образует взаимно ортогональную систему, а направления p^0, p^1, \dots, p^{N-1} Γ -сопряжены.

Доказательство. Пользуясь соотношением (7.59), (7.62), (7.61), лемму можно доказать методом математической индукции [10].

Действительно, по построению, p^1 сопряжен с p^0 . Кроме того, $p^0 = -\nabla f^0$, а по теореме Лагранжа ∇f^1 ортогонально p^0 , следовательно, ∇f^1 ортогонально ∇f^0 . Таким образом, для двух векторов лемма верна.

Предположим, что при $k < (N-1)$ векторы в системе p^0, p^1, \dots, p^k взаимно сопряжены, а векторы $\nabla f^0, \nabla f^1, \dots, \nabla f^k$ — взаимно ортогональны. Покажем, что эти свойства сохраняются у данных систем векторов при включении в них p^{k+1} и ∇f^{k+1} .

Рассмотрим значения $i < k$, тогда

$$((\nabla f^{k+1})^T \nabla f^i) = ((\nabla f^k + x^k \Gamma p^k)^T \nabla f^i) = x^k (\Gamma p^k)^T (-p^i + \beta_{i-1} p^{i-1}) = 0.$$

Равенство нулю получается за счет сопряженности p^k с векторами p^i и p^{i-1} .

Рассмотрим теперь $i=k$. Аналогично предыдущему $((\nabla f^{k+1})^T \nabla f^k) = ((\nabla f^k + x^k \Gamma p^k)^T \nabla f^k) = 0$. Равенство нулю можно получить, используя выражение из (7.61) для величины x^k .

Осталось доказать сопряженность системы векторов p^i для $i=1, \dots, k+1$. Сопряженность двух последних векторов следует из способа их построения. Осталось рассмотреть только $i < k$.

$$((p^{k+1})^T \Gamma p^i) = (-\nabla f^{k+1} + \beta_k p^k)^T \Gamma p^i = (-\nabla f^{k+1})^T \Gamma p^i = (-\nabla f^{k+1})^T (\nabla f^{i+1} - \nabla f^i) / x^i = 0.$$


Последнее равенство нулю вытекает из уже доказанной ортогональности градиентов.

Метод сопряженных направлений для положительно определенной квадратичной формы $f(y)$ построен. Однако, в формулу (7.60) для вычисления коэффициента β^{k-1} вошла неизвестная матрица Γ . Это не является существенным, поскольку формула (7.60) может быть переписана в другом виде. Чтобы показать это, выразим Γp^{k-1} в числителе (7.60) из (7.62), а p^{k-1} в знаменателе (7.60) из (7.59). Тогда

$$\begin{aligned} \beta_{k-1} &= (\nabla f^k, (\nabla f^k - \nabla f^{k-1})) / (x^{k-1} (-\nabla f^{k-1} + \beta_{k-2} p^{k-2})^T \Gamma p^{k-1}) = \\ &= (\nabla f^k, \nabla f^k) / (-(\nabla f^{k-1})^T x^{k-1} \Gamma p^{k-1}). \end{aligned}$$

Выражая $x^{k-1} \Gamma p^{k-1}$ из (7.62) и пользуясь ортогональностью ∇f^k и ∇f^{k-1} , окончательно получим

$$\beta_{k-1} = \|\nabla f^k\|^2 / \|\nabla f^{k-1}\|^2 \quad (7.63)$$

 **Замечание.** Построенный метод определяет минимум любой квадратичной функции с положительно определенной матрицей Гессе за N шагов. Отметим, что для определения x^k должен быть использован "аккуратный" одномерный поиск (т.е. параметр η одномерного поиска должен быть выбран близким к нулю).

Применение метода сопряженных градиентов к достаточно произвольной функции $f(y)$, естественно, не может обеспечить конечность процедуры поиска

минимума. После выполнения серии из N шагов метод, как правило, повторно запускается из последней найденной точки. Соответствующий алгоритм может быть записан следующим образом.

АЛГОРИТМ МЕТОДА СОПРЯЖЕННЫХ ГРАДИЕНТОВ ФЛЕТЧЕРА–РИВСА.

ШАГ 0. Задаются $\varepsilon > 0$ — параметр останова, $0 < \mu < \eta < 1$, $0 < \sigma < 1$ — параметры одномерного поиска, y^0 — начальная точка.

ШАГ 1. Вычисляются $f^0 = f(y^0)$, $\nabla f^0 = \nabla f(y^0)$, $p^0 = -\nabla f^0$, $k = 0$.

ШАГ 2. Если $(\nabla f^k, p^k) \geq 0$, то направление p^k не является направлением локального убывания функции, поэтому заменяем $p^0 = -\nabla f^k$, $y^0 = y^k$ и полагаем $k = 0$. Иначе направление p^k сохраняем. Переходим на шаг 3.

ШАГ 3. Вычисляется величина коэффициента одномерного шага x^k методом "аккуратного" одномерного поиска. Определяются

$$y^{k+1} = y^k + x^k p^k$$
$$f^{k+1} = f(x^{k+1}), \nabla f^{k+1} = \nabla f(x^{k+1}).$$

Полагается $k = k + 1$.

ШАГ 4. Проверяется критерий останова: при $\|\nabla f^k\| \leq \varepsilon$ поиск прекращается и y^k выдается как оценка решения; при $\|\nabla f^k\| > \varepsilon$ переходим к шагу 5.

ШАГ 5. Если $k = N$, полагается $y^0 = y^N$ и происходит возврат на шаг 1.

Если $k < N$, то переходим на шаг 6.

ШАГ 6. Вычисляем β^{k-1} по формуле (7.63) и p^k по формуле (7.59). Переходим на шаг 2.

Что известно о скорости сходимости построенного метода? Можно показать [39], что для функций из класса $\Phi_{m,M}$, описанного в (7.22), метод Флетчера-Ривса сходится со сверхлинейной скоростью.




Замечание. Метод чувствителен к "аккуратности" одномерного поиска и нарушению положительной определенности матрицы вторых производных минимизируемой функции. В указанном случае метод может построить направление p^k , не являющееся направлением локального убывания функции. В этом случае p^k заменяется на антиградиентное направление. Учет этой ситуации происходит на шаге 2 описания алгоритма.

7.7. Некоторые методы прямого поиска для негладких задач

В отличие от рассмотренных ранее, методы прямого поиска не используют каких-либо предположений о гладкости минимизируемой функции. Она может не только не иметь производных, но может содержать разрывы. При поиске минимума эти методы измеряют только значения функции. Поскольку гладкости нет, то при выборе направлений смещения методы не могут использовать аппроксимаций функции по результатам ее измерения. Правила размещения измерений в них основываются на некоторых эвристических логических схемах.

Наиболее популярными в практике расчетов являются следующие методы прямого поиска: Хука-Дживса [46], метод деформируемого многогранника Нелдера-Мида [47] и его модификация – комплексный метод Бокса. Нужно заметить, что последний метод применим только к выпуклым функциям. Поэтому здесь он не рассматривается. Ниже будут описаны первые два метода.

 **Замечание.** Несмотря на кажущуюся простоту и теоретическую необоснованность методов прямого поиска, они хорошо зарекомендовали себя в реальных расчетах.

Это можно объяснить следующим образом. Многие методы гладкой оптимизации чрезвычайно чувствительны к наличию вычислительных ошибок в значениях функций, превращающих теоретически гладкую функцию в фактически негладкую. За счет этого в реальных расчетах они зачастую утрачивают те положительные свойства, которые для них обещает теория. Использование методов прямого поиска позволяет в этих условиях добиться лучших результатов.

7.7.1. Метод Нелдера–Мида

В методе Нелдера–Мида вокруг начальной точки поиска в пространстве переменных размещается начальный симплекс – конфигурация из $(N+1)$ -й точки (в пространстве R^2 они образуют вершины треугольника, а в R^3 – вершины пирамиды). Затем происходит перемещение симплекса путем отражения вершины с наибольшим значением функции относительно центра тяжести противоположащего основания симплекса. При этом используются специальные операции, связанные с растяжением симплекса в направлении убывания функции и операции сжатия при неудачных пробных перемещениях. Дадим формальное описание алгоритма.

АЛГОРИТМ МЕТОДА НЕЛДЕРА–МИДА

ШАГ 0. Задаем векторы h^1, h^2, \dots, h^{N+1} , определяющие положение вершин стандартного симплекса с центром в начале координат, и числа S_1, S_2, \dots, S_{N+1} , определяющие размеры начального симплекса; $\varepsilon_y > 0$, $\varepsilon_f > 0$ – параметры останова; a, b, c, d – параметры отражения, растяжения, сжатия к основанию, сжатия к лучшей вершине ($a > 0$, $b > 1$, $0 < c < 1$, $0 < d < 1$). Задаем также начальную точку y^0 .

ШАГ 1. Формируем начальный симплекс с координатами вершин y^1, \dots, y^{N+1}

$$y^j = y^0 + S_j h^j; \quad (j = 1, \dots, N+1).$$

Вычисляем $f^j = f(y^j)$. (При этом в y^0 вычисление не выполняется).

ШАГ 2. Определяем номера худшей и лучшей вершины

$$f^h = \max\{f_j : j=1, \dots, N+1\}; \quad f^e = \min\{f_j : j=1, \dots, N+1\}.$$

ШАГ 3. Определяем центр тяжести основания

$$\bar{y} = \frac{1}{N} \left(\sum_{j=1, j \neq h}^{N+1} y^j \right).$$

ШАГ 4. Проверяем критерий останова. Вычисляем

$$\bar{y} = \frac{1}{N+1} \left(\sum_{j=1}^{N+1} y^j \right), \quad \bar{f} = \frac{1}{N+1} \left(\sum_{j=1}^{N+1} f^j \right), \quad \delta_y = \frac{1}{N+1} \left(\sum_{j=1}^{N+1} (y^j - \bar{y})^2 \right)^{1/2},$$

$$\delta_f = \frac{1}{N+1} \left(\sum_{j=1}^{N+1} (f^j - \bar{f})^2 \right)^{1/2}.$$

Если $\delta_y < \varepsilon_y$ и $\delta_f < \varepsilon_f$, то выполняем останов, выдаем оценку решения y^e, f^e .
 Если условия останова не выполнены, переходим на шаг 5.

ШАГ 5. Выполняем отражение с коэффициентом $a > 0$

$$y^* = \bar{y} + a(y^h - \bar{y})$$

и вычисляем $f^* = f(y^*)$.

ШАГ 6. Если $f^* > f^e$, то переходим на шаг 7. Если $f^* \leq f^e$, то выполняем растяжение

$$y^{**} = \bar{y} + b(y^* - \bar{y}), \quad b > 1, \quad f^{**} = f(y^{**});$$

при $f^{**} \leq f^*$ заменяем $y^h := y^{**}, f^h := f^{**}$ и переходим на шаг 2;
 при $f^{**} > f^*$ заменяем $y^h := y^*, f^h := f^*$ и переходим на шаг 2.

ШАГ 7. Если для любого $j = 1, \dots, N+1$, но $j \neq e$, выполняется $f^e < f^* < f^j$, то заменяем $y^h := y^*, f^h := f^*$ и переходим на шаг 2, иначе — на шаг 8.

ШАГ 8. Если $f^* < f^h$, то выполняем сжатие к основанию. Для этого вычисляем

$$y^\wedge = \bar{y} + c(y^h - \bar{y}), \quad f^\wedge = f(y^\wedge), \quad 0 < c < 1,$$

заменяем $y^h := y^\wedge, f^h := f^\wedge$ и переходим на шаг 2.

Если $f^* \geq f^h$, то выполняем сжатие к лучшей вершине:

$$y^j := y^e + d(y^j - y^e), \quad 0 < d < 1, \quad f^j := f(y^j) \quad (j = 1, \dots, N+1), \quad j \neq e$$

Переходим на шаг 2.

Авторы метода рекомендовали следующие значения параметров $a = 1$; $b = 1,5$; $c = 0,5$; $d = 0,5$ (кстати, метод чувствителен к их изменениям).

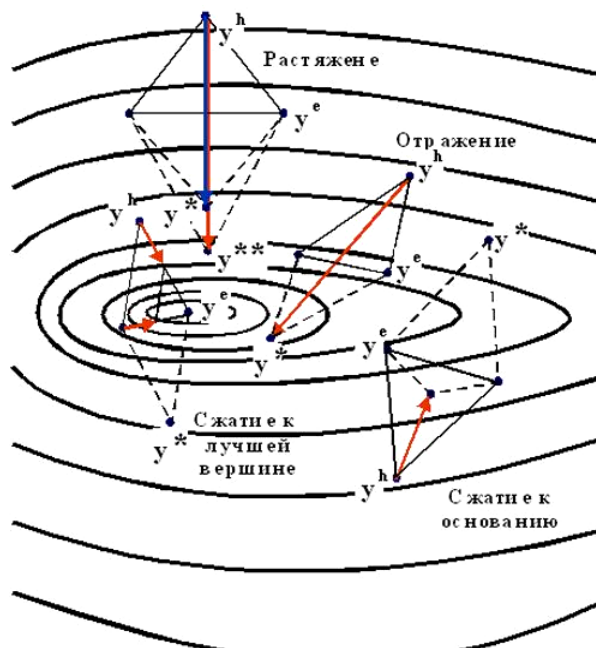


Рис. 7.17. Типовые операции с симплексом в методе Нелдера–Мида

Преобразования симплекса в пространстве R^2 при операциях отражения, растяжения и сжатия показаны на рис.7.17.

Метод Нелдера–Мида имеет тот недостаток, что для сильно овражных функций может происходить вырождение симплекса, особенно при числе переменных $N > 2$.

Термин «вырождение» означает, что все точки симплекса с некоторого шага размещаются в многообразии размерности меньшей, чем N , или же попадают в малую его окрестность, величина которой много меньше расстояния между точками симплекса.

7.7.2. Метод Хука-Дживса

В этом разделе приводится краткое описание метода Хука-Дживса, который был специально разработан именно для задач с оврагами [46]. В этом методе поиск минимума на каждом шаге происходит в результате смещения вдоль некоторого направления – образца (*шаг по образцу*), которое строится, а затем корректируется в результате специальных пробных покоординатных перемещений, называемых построением конфигурации.

Построение конфигурации из точки z осуществляет отображение z в точку $\bar{y} = F(z)$, где F – оператор построения конфигурации. Он устроен так, что направление $(\bar{y} - z)$ является направлением убывания функции f в окрестности z . Для описания оператора F введем следующие обозначения: e^i – i -й координатный орт, h – параметр, определяющий величину координатного перемещения. Тогда переход от z к y осуществляется согласно следующему алгоритму.

АЛГОРИТМ ПОСТРОЕНИЯ КОНФИГУРАЦИИ $\bar{y} = F(z)$:

ШАГ 0. Полагаем $\bar{y} = z$.

ШАГ 1. Для i от 1 до N выполнить:

если $f(\bar{y} + he^i) < f(\bar{y})$, то полагаем $\bar{y} := \bar{y} + he^i$, иначе, если $f(\bar{y} - he^i) < f(\bar{y})$, то $y := \bar{y} - he^i$.

На рис.7.18 показаны примеры построения конфигураций для нескольких случаев положения точки z . На рисунке пунктирными линиями отмечены пробные перемещения, не приведшие к уменьшению значения функции. Приведем пошаговое описание метода.

АЛГОРИТМ МЕТОДА ХУКА-ДЖИВСА:

ШАГ 0. Задаются начальная точка y^0 , параметр останова $\varepsilon > 0$ параметр построения конфигурации $h \gg \varepsilon$, а также параметр увеличения шага $\alpha = 2$.

ШАГ 1. Полагаем $z^1 = y^0$, $k = 0$.

ШАГ 2. Строим конфигурацию $y^{k+1} = F(z^{k+1})$.

ШАГ 3. Если $f(y^{k+1}) < f(y^k)$, то $k := k+1$ и переходим на шаг 4, иначе, если $h \leq \varepsilon$, выполняем ОСТАНОВ поиска, если $h > \varepsilon$, то дальнейшие действия зависят от того, как была построена точка y^{k+1} : строилась ли конфигурация с использованием шага по образцу (в этом случае $k > 0$) или она строилась от точки y^0 (в этом случае $k = 0$). Если окажется, что $k=0$, то сокращаем h вдвое ($h := h/2$) и переходим на шаг 1, если же $k > 0$, то полагаем $y^0 = y^k$, $k=0$ и также переходим на шаг 1.

ШАГ 4. Выполняем шаг по образцу $z^{k+1} = y^k + \alpha(y^k - y^{k-1})$ и переходим на шаг 2.

Одна из возможных ситуаций, связанных с использованием этого метода, показана на рис.7.18.

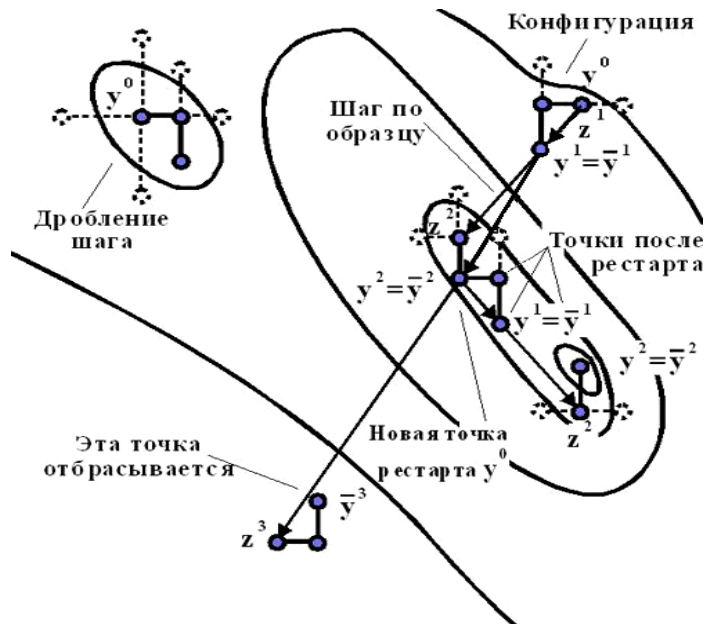


Рис. 7.18. Перемещения точки в методе Хука–Дживса из двух начальных положений y^0

Можно следующим образом пояснить смысл действий, выполняемых на шагах 2,3,4. Шаг 4 введен для того, чтобы метод обладал способностью быстрого увеличения величины смещения на одной итерации в том случае, когда точка y^k находится достаточно далеко от решения. При этом, прежде чем сделать z^{k+1} текущей точкой итерации, из точки z^{k+1} выполняется построение конфигурации.

За счет этого, в случае получения точки с $f(y^{k+1}) < f(y^k)$, следующий шаг по образцу в общем случае будет выполняться в измененном, по отношению к предыдущему, направлении, что позволяет методу адаптироваться к изменению рельефа функции.

Наконец, при получении значения $f(y^{k+1}) \geq f(y^k)$ на шаге 3 совершается попытка запустить метод сначала из точки $y^0 = y^k$ – лучшей найденной точки. При этом, если такой попытки еще не было, то параметр h не изменяется, а если она уже была, h предварительно сокращается вдвое.

7.8. Специальные методы учета линейных ограничений в гладких задачах локальной оптимизации

В предыдущих разделах были подробно рассмотрены несколько групп методов поиска локально-оптимальных решений в задачах без ограничений. В действительности, ограничения почти всегда присутствуют. В первую очередь это относится к ограничениям, определяющим диапазоны изменения переменных. Если задача имеет естественнонаучную или техническую природу, то ограничения на переменные возникают из их «физического» смысла, не позволяя переменным принимать сколь угодно большие или сколь угодно малые значения.

Кроме двусторонних ограничений на переменные вида (7.3) в задачах могут присутствовать функциональные ограничения общего или специального вида (7.2). В дальнейшем мы будем разделять два вида этих ограничений.

Общие ограничения, определяемые через функции ограничений равенств и неравенств, проще всего учесть с помощью одного из рассмотренных в главе 3 общих методов учета ограничений, сводящего задачу с функциональными ограничениями к серии задач без функциональных ограничений (например, метод внешнего штрафа [5, 8] — см. раздел 3.2, метод модифицированных функций

Лагранжа [2] — см. раздел 3.3 и другие [5, 23] — см. раздел 3.4). Такой подход позволяет применить к решению задач с ограничениями методы, ранее разработанные для задач без ограничений. Однако, если ограничения имеют специальный вид, допускающий их явный учет за счет соответствующей модификации применяемых методов, то способы их учета целесообразно изучить отдельно.

К специальным ограничениям в первую очередь следует отнести произвольные линейные ограничения (равенства и неравенства), а также — двухсторонние ограничения на переменные.

7.8.1. Специальные методы учета линейных равенств

Рассмотрим задачу только с линейными ограничениями–равенствами

$$\begin{aligned} f(y) \rightarrow \min, \quad y \in Y, \\ Y = \{y \in R^N : (h_i, y) = c_i, \quad (i = 1, \dots, p)\}. \end{aligned} \quad (7.64)$$

Будем считать, что $Y \neq \emptyset$ и допустимое множество Y является многообразием положительной размерности, а линейно зависимые ограничения из постановки задачи исключены.

Если ввести матрицу H , столбцами которой являются векторы h_1, \dots, h_p , то ограничения запишутся в форме $H^T y = c$. Пусть R_H — подпространство, натянутое на набор векторов $\{h_1, \dots, h_p\}$, а R_W — его ортогональное дополнение. Очевидно, что матрица, столбцы которой образуют базис в R_H , может быть выбрана равной H . Базисную матрицу того же типа для R_W обозначим через W . Таким образом, всегда $H^T W = 0$.

Лемма 7.4. Пусть $H^T = (V; Z)$, где V — невырожденная квадратная матрица (возможно, такое представление потребует перенумерации переменных), тогда можно построить W в виде блочной матрицы следующего вида (E — единичная матрица) $W = \begin{pmatrix} -V^{-1}Z \\ E \end{pmatrix}$.

Доказательство очевидно, т.к. непосредственным перемножением блочных матриц получим $H^T W = -Z + Z = 0$. Лемма доказана.

Представим y в виде разложения

$$y = H y_H + W y_W. \quad (7.65)$$

Используя $H^T y = c$, видим, что условием допустимости точки y является использование в (7.65) значения

$$y_H = y_H^* = (H^T H)^{-1} c. \quad (7.66)$$

В результате исходная задача (7.64) сводится к задаче без ограничений в пространстве размерности $N-p$:

$$\tilde{f}(y_w) \rightarrow \min, \quad y_w \in R^{N-p}, \quad \tilde{f}(y_w) = f(H y_H^* + W y_w). \quad (7.67)$$

Рассмотрим особенности применения методов гладкой оптимизации в задаче (7.64) с учетом возможности ее сведения к форме (7.67). В первую очередь необходимо использовать очевидную лемму.

Лемма 7.5. Существует следующая связь между градиентами и матрицами Гессе в пространстве исходных переменных (в точке $y^k = \mathbf{H}y_H^* + \mathbf{W}y_w^k$) и новых переменных в точке y_w^k

$$\nabla \tilde{f}^k = \nabla \tilde{f}(y_w^k) = \mathbf{W}^T \nabla f(y^k), \quad \Gamma_k^{\tilde{f}} = \Gamma^{\tilde{f}}(y_w^k) = \mathbf{W}^T \Gamma^f(y^k) \mathbf{W}.$$

Пусть далее в процессе применения одного из методов безусловной локальной оптимизации в пространстве переменных y_w было выбрано направление d_w^k , тогда направлением поиска в исходном пространстве будет

$$d^k = \mathbf{W} d_w^k, \quad (7.68)$$

и величина смещения будет определяться в исходном пространстве для направления d^k применительно к задаче (7.64).

В зависимости от метода гладкой оптимизации, выбранного для решения задачи (7.67), выбор направления выполняется по следующим правилам.

А. В методе Ньютона–Рафсона

$$d_w^k = \left(\Gamma_k^{\tilde{f}} \right)^{-1} \left(-\nabla \tilde{f}^k \right).$$

В. В методе Ньютона–Рафсона с модификацией матрицы Гессе строится модифицированное разложение Холесского для $\Gamma_k^{\tilde{f}} : \Gamma_k^{\tilde{f}} \Rightarrow \bar{L}_k^{\tilde{f}} \bar{D}_k^{\tilde{f}} \left(\bar{L}_k^{\tilde{f}} \right)^T$ и определяется d_w^k из последовательного решения двух систем с треугольными матрицами

$$\bar{L}_k^{\tilde{f}} v = -\nabla \tilde{f}^k, \quad \left(\bar{D}_k^{\tilde{f}} \left(\bar{L}_k^{\tilde{f}} \right)^T \right) d_w^k = v.$$

С. В квазиньютоновских методах матричные оценки размера $(N-p) \times (N-p)$ вычисляются в пространстве R_w , а именно, на общем шаге определяются приращения

$$z_w^k = \nabla \tilde{f}^{k+1} - \nabla \tilde{f}^k, \quad \Delta_w^k = y_w^{k+1} - y_w^k = x^k d_w^k$$

и выполняется коррекция матричных оценок

$$G_{k+1}^{\tilde{f}} = G_k^{\tilde{f}} + U(z_w^k, \Delta_w^k)$$

по одной из формул (*B*, *DFP* или *BFGH*) пункта 7.5.1. Далее выбирается


$$d_w^{k+1} = \left(G_{k+1}^{\tilde{f}} \right)^{-1} \left(-\nabla \tilde{f}^{k+1} \right).$$

Д. В квазиньютоновских методах с модификацией матричных оценок перед определением d_w^{k+1} выполняется модифицированное преобразование Холесского для матриц $G_{k+1}^{\tilde{f}}$.

Е. В методе Флетчера–Ривса выполняется построение сопряженных направлений для пространства R_w

$$p_w^{k+1} = -\nabla \tilde{f}^{k+1} + \beta_k p_w^k, \quad \beta_k = \frac{\|\nabla \tilde{f}^{k+1}\|^2}{\|\nabla \tilde{f}^k\|^2}$$

и выбирается $d_w^{k+1} = p_w^{k+1}$.

 **Замечание.** Поскольку методы из группы квазиньютоновских (переменной метрики) и метод Флетчера–Ривса выполняют циклические рестарты через число шагов, кратное размерности пространства, то при реализации методов (С)–(Е) необходимо выполнять рестарты на каждом $(N-p+1)$ –м шаге.

7.8.2. Специальные методы учета линейных неравенств, методы активного набора

Рассмотрим задачи только с линейными неравенствами

$$\begin{aligned} f(y) \rightarrow \min, \quad y \in Y, \\ Y = \{y \in R^N : (g_j, y) \leq g_j^+, \quad (j = 1, \dots, m)\}. \end{aligned} \quad (7.69)$$

Если ввести матрицу G , столбцами которой являются векторы g_1, \dots, g_m , то ограничения запишутся в форме $G^T y \leq g^+$.

Пусть y^* – решение этой задачи. Множество номеров ограничений–неравенств, активных в точке решения, обозначим через $J^* = J(y^*) = \{j_1, \dots, j_{r^*}\}$. Остальные ограничения не влияют на решение, поэтому вместо задачи (7.69) можно решать эквивалентную задачу с линейными ограничениями–равенствами, подобную (7.64), где принято

$$p=r^*, \quad h_i = g_{j_i}, \quad c_i = g_{j_i}^+ \quad (i=1, \dots, r^*).$$

Таким образом, для ее решения можно применить рассмотренный в пункте 7.8.1 подход, однако сложность состоит в том, что набор индексов J^* заранее не известен.

Для преодоления подобных трудностей обычно применяют *метод рабочего набора*. Его использование тесно связано с условиями экстремума первого порядка, рассмотренными в пункте 2.2.1 (см. теоремы.2.5–2.7, которые нужно интерпретировать применительно к рассматриваемой однокритериальной задаче). Ниже приводится краткое описание этого метода.

ОБЩАЯ СХЕМА МЕТОДА РАБОЧЕГО НАБОРА

ШАГ 0. В качестве гипотезы принимается некоторый рабочий набор $J = \{j_1, \dots, j_r\} \subseteq \{1, \dots, m\}$ и начальная точка y^0 , что $J(y^0) = J$ (вначале можно принять $r=0, J = \emptyset$).

ШАГ 1. Номер шага k полагается равным 0, строится вспомогательная задача с линейными ограничениями–равенствами

$$f(y) \rightarrow \min, \quad H_J^T y = c_J, \quad (7.70)$$

где $H_J = (g_{j_1}, \dots, g_{j_r}), c_J = (g_{j_1}^+, \dots, g_{j_r}^+)^T$.

ШАГ 2. Ищется решение задачи (7.70) начиная с точки y^0 с помощью выбранного метода локальной оптимизации. При выполнении методом каждого шага, в исходной задаче (7.69) проводится контроль нарушения ограничений–неравенств с номерами $j \notin J$.

Если в процессе очередного шага точка поиска y^k_J выходит на границу ограничения с номером $j \notin J$ и дальнейшее смещение приводит к его нарушению, то полагается $y^0_J = y^k_J$, выполняется расширение рабочего набора

$$J := J \cup \{j\}$$

и производится возврат к шагу 1.

Если найдено решение y^*_J задачи (7.70), то выполняется переход на шаг 3.

ШАГ 3. На основе условий Куна–Таккера выполняется оценка множителей Лагранжа. Для этого записывается линейная система относительно λ следующего вида


$$-\nabla f(y_J^*) = H_J \lambda \quad (7.71)$$

Если покомпонентно $\forall i \in \{1, \dots, r\} \lambda_i \geq 0$, то выполняется останов процесса вычислений и точка y_J^* выдается в качестве решения задачи.

Если $\exists i \in \{1, \dots, r\} \lambda_i < 0$, то согласно замечанию к теореме 2.5, необходимо вывести процесс поиска с границ ограничений с соответствующими номерами для возможности дальнейшего уменьшения значений функции в допустимой области. Для этого полагается $y_J^0 = y_J^*$, уменьшается число элементов в рабочем наборе

$$J := J \setminus \{j_i \in J : \lambda_i < 0\}$$

и производится возврат к шагу 1.

 **Замечание.** За счет погрешностей численного определения решений y_J^* вспомогательных задач, система уравнений (7.71) будет являться переопределенной. Перед ее решением обычно исключают часть уравнений, оставляя вместо матрицы H_J квадратный блок полного ранга размером $r \times r$.

Проанализируем величину ошибки в определении истинного значения λ^* множителей Лагранжа. Пусть $\Delta\lambda = \lambda^* - \lambda$, а ε – ошибка решения вспомогательных задач. Тогда

$$-\nabla f(y_J^*) + O(\|\varepsilon\|) = H_J \lambda^* + H_J \Delta\lambda,$$

т.е.

$$H_J \Delta\lambda = O(\|\varepsilon\|).$$

Следовательно, погрешность $\Delta\lambda$ является величиной того же порядка, что и ε , но на ее величину дополнительно влияет степень обусловленности матрицы H_J .

7.8.3. Особенности применения методов локального поиска при двусторонних ограничениях на переменные

Двусторонние ограничения на переменные являются частным случаем линейных неравенств, и их можно было бы учесть на основе общей методики для ограничений такого типа. Однако их вид настолько прост, с одной стороны, а, с другой стороны, специфичен, что наиболее правильным решением является отдельное описание способа работы с ними. Наличие таких ограничений неизбежно приводит к пересмотру ранее построенных в разделах 7.2–7.6 алгоритмов и созданию их модификаций для задач с двусторонними ограничениями.

В зависимости от типа метода его модификация выполняется по-разному. Общие принципы построения модифицированных методов можно предложить для методов гладкой оптимизации, а для методов прямого поиска приходится использовать в каждом случае свои уникальные подходы.

7.8.3.1 Особенности учета двусторонних ограничений на переменные в методах гладкой оптимизации

Характерными чертами многих методов гладкой оптимизации (для задач без ограничений специального вида) являются:

- выполнение рестартов из последней достигнутой точки через определенное число шагов, которое зависит от размерности пространства поиска (например, методы квазиньютоновского типа, метод сопряженных градиентов — см. материал пункта 7.8.1);
- выполнение одномерного поиска в выбранном направлении (в большинстве методов) либо перемещения в выбранном направлении с заранее заданным коэффициентом величины шага (метод Ньютона).

Рассмотрим теперь задачу с двусторонними ограничениями на переменные

$$f(y) \rightarrow \min, y \in D, \quad D = \{y \in R^N : a \leq y \leq b\}. \quad (7.72)$$

Наличие таких ограничений будет оказывать влияние на реализацию каждого из этих двух процессов. А именно, процессы одномерных перемещений могут выводить на фрагменты границы области D . После выхода на границу поиск должен продолжаться на линейном координатном многообразии меньшей размерности. Одномерные перемещения в этом многообразии могут выводить процесс поиска на ограничения по другим переменным, что будет приводить к дальнейшему понижению размерности многообразия поиска. Кроме того, поиск на возникающих многообразиях, кроме контроля пересечения границ области, будет иметь также ту особенность, что методы, периодически выполняющие рестарты, должны будут производить их чаще, чем при поиске во всем пространстве. Это связано с тем, что количество шагов между рестартами определяется размерностью многообразия, на котором выполняется поиск. Еще одним важным моментом является правило возврата процесса поиска с текущего многообразия на многообразии (или в пространство) более высокой размерности в том случае, когда это приводит к уменьшению значения функции. Все эти особенности уже присутствовали в методе рабочего набора, рассмотренном в предыдущем разделе, однако их реализация при учете двусторонних ограничений может быть выполнена в более простой форме.

Дадим описание правил выполнения следующих операций:

- учет выходов на новые фрагменты границы при одномерных перемещениях;
- организация поиска на многообразиях размерности $n < N$;
- определение моментов возврата с многообразий текущей размерности n в многообразия большей размерности.

Для описания текущего многообразия поиска введем два множества J_a и J_b . В последующем они будут содержать наборы номеров переменных, по которым текущая точка поиска выведена на нижние или верхние граничные значения. Если хотя бы одно из этих множеств не пусто, то их совокупность идентифицирует линейное многообразие размерности $n < N$, в котором происходит поиск. Это многообразие соответствует фиксации части компонент y_i вектора y на граничных значениях. Перед началом поиска множества J_a и J_b должны быть пусты.

Пусть в результате очередного шага, выполненного в направлении d^{k-1} , процесс поиска вышел на границу области D в точке y^k . Пусть этот участок границы имеет размерность $n < N$. Для текущей точки y^k и направления d^{k-1} скорректируем множества J_a и J_b , включив в них номера переменных по которым процесс поиска вышел, соответственно, на верхние или нижние границы их изменения. Формально определим правило коррекции следующим образом.

$$J_a := J_a \cup \{i: 1 \leq i \leq N; y_i^k = a_i; d_i^{k-1} < 0\} \quad (7.73)$$

$$J_b := J_b \cup \{i: 1 \leq i \leq N; y_i^k = b_i; d_i^{k-1} > 0\} \quad (7.74)$$

Первое изменение множеств J_a и J_b соответствует переходу процесса поиска из пространства размерности N на линейное многообразие меньшей размерности за счет фиксации дополнительных компонент y_i вектора y . Введем базис в пространстве свободных (нефиксированных) координат. Для этого из набора единичных координатных ортов e^1, \dots, e^N выделим те векторы e^j для которых $j \notin (J_a \cup J_b)$. Составим из этих векторов матрицу W_k , используя их как вектор-столбцы

$$W_k = (e^{j^1}, \dots, e^{j^n}).$$

Введем вектор переменных y_w для возникшего линейного подпространства $R_w = R^n$. Переход процесса поиска на линейное многообразие из исходного пространства R^N равносильна замене переменных $y = y^k + W_k y_w$ с переходом к поиску в пространстве переменных y_w .

Вычисляя в старых переменных значения ∇f^k и Γ_k легко пересчитать их в соответствующие значения в новых переменных, используя известные соотношения

$$\nabla_{y_w} f^k = (W_k)^T \nabla f^k; \quad \Gamma_{y_w k} = (W_k)^T \Gamma_k W_k$$

Нужно обратить внимание на то, что в пространстве новых переменных методы гладкой оптимизации должны быть запущены заново из точки $y_w^0 = 0$, соответствующей текущей точке поиска y^k . Заметим также, что если используемый метод включает рестарты, то при поиске минимума по переменным y_w эти рестарты необходимо выполнять через число шагов, согласованное с размерностью n многообразия поиска. Например, для квазиньютоновских методов и метода сопряженных градиентов рестарты следует выполнять через n шагов.

Рассмотрим процесс поиска на многообразии. Выбор очередного направления поиска в переменных y_w происходит по обычным правилам, характерным для выбранного метода. Однако, при реализации этих правил необходимо все векторы и матрицы использовать для размерности n нового пространства, т.е., в частности, вместо значений ∇f^k и Γ_k необходимо использовать $\nabla_{y_w} f^k$ и $\Gamma_{y_w k}$.

Работа методов в пространстве переменных y_w имеет дополнительную специфику, связанную с тем, что в действительности метод решает исходную N -мерную задачу с двусторонними ограничениями, наличие которых влияет на правила выполнения шага методом. Допустим, метод находится в точке y^k и, согласно правилам применяемого алгоритма, в пространстве переменных y_w выбрано направление поиска d_w^k . Тогда (также как при учете произвольных линейных ограничений) выполняется пересчет направления d_w^k в направление d^k в исходном пространстве переменных:

$$d^k = W_k d_w^k.$$

После выбора направления поиска происходит перемещение в этом направлении. Смещение выполняется в многообразии, соответствующем множествам J_a и J_b . В зависимости от типа метода это перемещение выполняется либо с фиксированным коэффициентом одномерного шага $x = const$, либо за счет поиска минимума вдоль выбранного направления. В обоих случаях учитываются ограничения на переменные.

Процедура одномерного поиска, приведенная в пункте 7.2.4, учитывает их автоматически. Если при ее выполнении точка $y^{k+1} = y^k + x^k d^k$, переместившись вдоль многообразия, выходит на границу области по новым переменным, то их

номера следует добавить в множества J_a или J_b , соответственно, дополнительно применив правила их коррекции (7.73), (7.74) для $k=k+1$.

Если же перемещение точки должно быть выполнено с фиксированным коэффициентом длины шага (как это происходит, например, в методе Ньютона), то, в случае выхода точки за пределы изменения переменных, коэффициент α длины шага уменьшается таким образом, чтобы точка $y^{k+1}=y^k+\alpha d^k$ оказалась на границе области D . Далее выполняется описанная выше коррекция множеств J_a и J_b .

Осталось рассмотреть случай, когда одномерное перемещение в пространстве переменных y_w не приводит к выходу точки y^{k+1} на новые фрагменты границы. В этом случае необходимо проверить условия возврата к многообразию или пространству более высокой размерности, который может быть выполнен за счет исключения из множеств J_a или J_b номеров части переменных. Правила исключения следующие:

$$J_a := J_a \setminus \{i \in J_a: \partial f(y^k)/\partial y_i < 0\} \quad (7.75)$$

$$J_b := J_b \setminus \{i \in J_b: \partial f(y^k)/\partial y_i > 0\}. \quad (7.76)$$

Условия исключения переменной в (7.75), (7.76) определяются тем, что в текущей точке поиска становится положительной проекция антиградиента функции f на внутреннюю (по отношению к области D) нормаль к гиперплоскости $y_i = a_i$. При выполнении этого условия существует направление смещения с этой гиперплоскости внутрь области, при котором функция f будет локально убывать. Таким образом, если в результате коррекции (7.75), (7.76) хотя бы одно из множеств J_a или J_b изменится, необходимо перейти к поиску в многообразии более высокой размерности, выполнив в нем рестарт метода из последней точки предшествующего поиска, аналогично тому, как это было описано выше.

7.8.3.2. Учет двусторонних ограничений в методах прямого поиска

В методах прямого поиска способ учета двусторонних ограничений уникален для каждого из этих методов. Более того, некоторые из них не имеют точных модификаций для задач с двусторонними ограничениями. Типичным примером является метод Нелдера–Мида. Некоторые же методы, например метод Хука–Дживса, напротив, легко обобщаются на такие задачи.

Рассмотрим принципы модификации для метода Хука–Дживса. В нем выполняются действия только двух типов: шаг по образцу и построение конфигурации.

С учетом ограничений, шаг по образцу выполняется таким образом, что при выходе рабочей точки за границы области D величина последнего смещения корректируется так, чтобы точка оказалась на границе D . При построении конфигурации в обычном методе координатные перемещения выполняются с шагом h . В модифицированном методе величины координатных перемещений не превосходят h , а в случае, если эти перемещения выводят из области D , заменяются на перемещения до границ этой области.

Рассмотрим метод Нелдера–Мида. На первый взгляд правила метода допускают аналогичный способ модификации. Однако при этом границы области будут трансформировать правила отражения и растяжения симплекса. Очевидно, что это будет способствовать быстрому его вырождению в окрестности границ области. Следовательно, такой подход не применим, хотя возможны специальные модификации этого метода для случая ограничений, рассмотренные, например, в

книге Ф. Гилла и У. Мюррея [13]. Кроме того, учет двусторонних ограничений на переменные в этом методе можно выполнить с использованием общего метода внешнего штрафа, который был рассмотрен в разделе 3.2.

Заключение

Те из читателей, которым хватило сил и желания освоить весь материал книги, легко обнаружат, что к моменту окончания последней главы ими оказались изучены все основные постановки задач математического программирования и многокритериальной оптимизации, ключевые факты из теории условий экстремума, распространенной на многокритериальный случай, теории двойственности, а также основные семейства вычислительных методов решения рассмотренных классов задач, включая методы учета ограничений различного вида, методы локальной, многоэкстремальной и многокритериальной оптимизации.

В материалах книги внимательный читатель наверняка обнаружил много малоизвестных результатов и идей, касающихся одношагово–оптимальных подходов к конструированию алгоритмов, методов редукции размерности, теории сходимости и анализа плотности размещения испытаний в методах многоэкстремальной оптимизации.

Авторы выражают надежду, что книга озалась полезной студентам, аспирантам и специалистам, интересующимся как методами так и теорией конечномерной оптимизации.

Лист регистрации изменений

Дата	Автор	Комментарии
08.11.01	Городецкий С.Ю.	Создание документа и раздела 7.7
15.08.01	Городецкий С.Ю.	Создание раздела 7.2
28.08.01	Городецкий С.Ю.	Создание раздела 7.3
24.07.02	Городецкий С.Ю.	Создание раздела 7.1
24.07.02	Городецкий С.Ю.	Внесение изменений в раздел 7.1
25.07.02	Городецкий С.Ю.	Внесение изменений в раздел 7.2
26.07.02	Городецкий С.Ю.	Внесение изменений в раздел 7.3
27.07.02	Городецкий С.Ю.	Внесение изменений в разделы 7.4–7.7
28.07.02	Городецкий С.Ю.	Внесение изменений в раздел 7.8
29.07.02	Городецкий С.Ю.	Внесение изменений в раздел 7.8
19.11.03 – 22.11.03	Городецкий С.Ю.	Добавление и оптимизация рисунков в разделах 7.1–7.6, корректурa этих разделов
23.11.03 – 26.11.03	Городецкий С.Ю.	Корректурa разделов 7.7–7.8

Литература

Основная литература

10. Базара М., Шетти К. Нелинейное программирование. Теория и алгоритмы. – М.: Мир, 1982.
11. Бертсекас Д. Условная оптимизация и методы множителей Лагранжа.– М.: Радио и связь, 1987
12. Васильев Ф.П. Численные методы решения экстремальных задач. – М.:Наука, 1982.
13. Гилл Ф., Мюррей У., Райт М. Практическая оптимизация.– М.: Мир, 1985.
14. Карманов В.Г. Математическое программирование. – М.:Наука, 1986.
15. Подиновский В.В, Ногин В.Д. Парето–оптимальные решения многокритериальных задач.– М.: Наука, 1982
16. Стронгин Р.Г. Численные методы многоэкстремальной оптимизации (информационно–статистические алгоритмы) – М.:Наука, 1978.
17. Сухарев А.Г., Тимохов А.В., Федоров В.В. Курс методов оптимизации. – М.:Наука,1986.
18. Strongin R.G.,Sergeev Ya.D. Global Optimization with Non–Convex Constraints. Sequential and Parallel Algorithms.– Dordrecht: Kluwer Academic Publishers. The Netherlands, 2000.

Дополнительная литература

19. Аоки М. Введение в методы оптимизации.– М.:Наука, 1977.
20. Батищев Д.И. Поисковые методы оптимального проектирования. М.: Советское радио, 1975.
21. Гермейер Ю.Б. Введение в теорию исследования операций. М.: Наука, 1971.
22. Гилл, у, Мюррей Численные методы условной оптимизации. — М.:Мир, 1977, п. 7.13.
23. Городецкий С.Ю., Неймарк Ю.И. О поисковых характеристиках алгоритма глобальной оптимизации с адаптивной стохастической моделью. // Проблемы случайного поиска./ Под ред. Л.А. Растригина. Рига: Зинатне, 1981. Вып.9 стр. 83-105.
24. Городецкий С.Ю. Сходимость и асимптотические оценки поведения для одного класса методов поиска. // Динамика систем. Динамика и управление. Горький: ГГУ, 1984.
25. Городецкий С.Ю. Метод поиска условного глобального минимума, основанный на простых вероятностных моделях минимизируемой функции и функций ограничений. //Динамика систем. Управление и оптимизация. Горький: ГГУ, 1987.
26. Городецкий С.Ю. Процедура решения задач многокритериальной оптимизации, основанная на простых адаптивных вероятностных моделях критериев. //Непараметрические и робастные методы в кибернетике и информатике. Материалы VII Всесоюзного семинара. Томск: изд-во ТГУ, 1990. Часть 1, с.156–159.
27. Городецкий С.Ю. Методы поиска множества решений многокритериальных задач на основе простых вероятностных моделей критериев. ННГУ, Н.Новгород, 1992. — Деп. в ВИНТИ 25.03.92, № 1024–В92, 18 с.
28. Городецкий С.Ю. Многоэкстремальная оптимизация на основе триангуляции области.// Математическое моделирование и оптимальное управление. Вестник Нижегородского государственного университета. Вып. 2(21). Н.Новгород: изд-во ННГУ, 1999, с.249–269.
29. Городецкий С.Ю. Методы многоэкстремальной оптимизации на основе триангуляции области поиска. //В кн.: Первая всероссийская научно-практическая конференция по вопросам решения научно-практических задач в промышленности ОПТИМ–2001. Сборник докладов. — С-Пб.: ЦНИИТС, 2001, с.191–196.
30. Гришагин В.А. Об условиях сходимости для одного класса алгоритмов глобального поиска. – В сб.: Тезисы докл. III Всес. семинара "Численные методы нелинейного программирования". Харьков:ХГУ, 1979, с. 82-84.

31. Гришагин В.А., Стронгин Р.Г. Оптимизация многоэкстремальных функций при монотонно унимодальных ограничениях. //Изв. АН СССР. Техническая кибернетика, 1984, №4, с. 203-208.
32. Евтушенко Ю.Г. Методы решения экстремальных задач и их применение в системах оптимизации. –М.:Наука, 1982.
33. Евтушенко Ю.Г., Потапов М.А. Методы численного решения многокритериальных задач. //ДАН, 1986, т.28, № 11.
34. Евтушенко Ю.Г., Потапов М.А. Численные методы решения многокритериальных задач. //В кн.: Кибернетика и вычислительная техника. Вып.3. — М.: Наука, 1987, с.209–218.
35. Евтушенко Ю.Г., Ратькин В.А. Метод половинных делений для глобальной оптимизации функций многих переменных. //Техническая кибернетика. 1987, №1, с.119–127.
36. Жиглявский А.А., Жилинскас А.Г. Методы поиска глобального экстремума.–М.: Наука, 1991
37. Жилинскас А.Г. Одношаговый байесовский метод поиска экстремума функций одной переменной. // Кибернетика 1975, № 1 с. 139–144.
38. Лузин Н.Н. Теория функций действительного переменного. — М.: Учпедгиз,1948.
39. Маркин Д.Л., Стронгин Р.Г. О равномерной оценке множества слабоэффективных точек в многоэкстремальных многокритериальных задачах оптимизации. //ЖВМиМФ. 1993,т.33,№ 2.
40. Методы поиска глобального экстремума. Методические указания./ Сост. Городецкий С.Ю.– Горький:, ГГУ, 1990.
41. Моцкус Й.Б. О байесовских методах поиска экстремума. //Автоматика и вычислительная техника, 1972, № 3, с.53–62.
42. Моцкус Й.Б. Исследование простого байесового алгоритма для решения многоэкстремальных задач. ИМиК АН ЛССР, Вильнюс, 1979. — Деп. В ВИНТИ 1979, № 4291–79, 8с.
43. Неймарк Ю.И. Динамические системы и управляемые процессы.— М.: Наука, 1978.
44. Неймарк Ю.И., Стронгин Р.Г. Информационный подход к задаче поиска экстремума функций. //Известия АН СССР. Техническая кибернетика, 1966, №1, с.17-26.
45. Немировский А.С., Юдин Д.Б. Сложность задач и эффективность методов оптимизации. – М.: Наука, 1979.
46. Пиявский С.А. Алгоритмы отыскания абсолютного минимума функций. //В кн.: Теория оптимальных решений. Вып.2. Киев, 1967, с.13-24.
47. Пиявский С.А. Один алгоритм отыскания абсолютного минимума функции. //ЖВМ и МФ, №4, 1972, с.888-896.
48. Пшеничный Б.Н., Данилин Ю.М. Численные методы в экстремальных задачах.– М.: Наука, 1975.
49. Сергеев Я.Д., Квасов Д.Е. Адаптивные диагональные кривые и их программная реализация. // Вестник ННГУ. Мат. Моделирование и оптимальное управление. Вып. 2(24), 2001, с.300–317.
50. Соболев И.М., Статников Р.Б. Выбор оптимальных параметров в задачах со многими критериями. — М.: Наука, 1981.
51. Стронгин Р.Г. Поиск глобального оптимума. –М.:Знание,1990.
52. Стронгин Р.Г., Гергель В.П., Городецкий С.Ю., Гришагин В.А., Маркина М.В. Современные методы принятия оптимальных решений. — Н. Новгород: Изд-во ННГУ, 2002.
53. Стронгин Р.Г., Маркин Д.Л. Минимизация многоэкстремальных функций при невыпуклых ограничениях.// Кибернетика №4, 1986, с.64-69.
54. Сухарев А.Г. Оптимальный поиск экстремума. – М.: Изд МГУ, 1975.
55. Уайлд Д.Дж. Методы поиска экстремума. 1967.
56. Химмельблау Д. Прикладное нелинейное программирование. –М.:Мир, 1975.
57. Черноусько Ф.Л., Меликян. Игровые задачи управления и поиска.–М.: Наука, 1978

58. Шор Н.З. О скорости сходимости метода обобщенного градиентного спуска с растяжением пространства. //Кибернетика, № 2, 1970.
59. Kushner H.J. A versatile Stochastic model of a function of unknown and time-varying form. //J. Math. anal and appl. v5, №1, 1962, pp.150-167.
60. Kushner H. A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise. //Transactions ASME, Ser. D, J. Basic Eng., 1964, vol.86, №.1, p.97-106.
61. Pinter J. Global optimization in Action. –Dordrecht: Kluwer Academic Publishers. The Netherlands, 1996.
62. Strongin R.G., Sergeyev Ya.D., Grishagin V.A. Parallel Characteristical Algorithms for Solving Problems of Global Optimization // Journal of Global Optimization,10, 1997, pp. 185-206.
63. Sergeyev Ya.D. On Convergence of "Divide the Best" Global Optimization Algorithms. //Optimization, Vol.44, 1998, pp.303-325.
64. Sergeyev Ya.D. An efficient strategy for adaptive partition of N-dimensional intervals in the framework of diagonal algorithms. // Journal of Optimization Theory and Applications, vol.107, №.1, 2000, pp. 145–168.
65. Sergeyev Ya.D., Grishagin V.A. A Parallel Method for Finding the Global Minimum of Univariate Functions. // Journal of Optimization Theory and Applications, vol.80, №.3, 1994.

Это будет указано в сносках

Митягин Б.С. Два неравенства для объемов выпуклых тел. – Мат. заметки, 1969, т.5, вып.1.

Белоглазов И.Н., Джанджагава Г.И., Чигин Г.П. Основы навигации по геофизическим полям /Под ред.А.А.Красовского.— М.: Наука, 1985.

