

Нижегородский государственный университет им. Н. И. Лобачевского

**Создание в Нижегородском государственном университете
межфакультетской магистратуры "Математические модели,
методы и программное обеспечение современных
компьютерных технологий"**

**Образовательный комплекс.
Современные методы принятия
оптимальных решений**

Электронный учебник

Авторский коллектив

Руководитель профессор, д.т.н. Гегель В.П.

Ответственный исполнитель доцент, к.ф.-м.н. Гришагин В.А.

Главный разработчик доцент, к.ф.-м.н. Городецкий С.Ю.

Разработчик доцент, к.ф.-м.н. Маркина

Оглавление

Глава 1. Математические модели принятия решений. Примеры прикладных задач	5
1.1. Методология принятия решений	5
1.2. Общая модель рационального выбора	6
1.3. Модели математического программирования и методы их анализа	9
1.4. Принципы создания учебно-исследовательских систем принятия решений	12
1.5. Прикладные задачи оптимального выбора	16
1.5.1. Задача определения местоположения и параметров движения объекта по измерениям высоты рельефа местности	16
1.5.2. Задача определения обобщенных координат манипулятора для заданного позиционирования схвата	18
Краткий обзор главы	24
Глава 2. Математические основы конструирования алгоритмов. Характеристические алгоритмы глобального поиска	25
2.1. Постановка задачи	25
2.2. Численные методы оптимизации	26
2.3. Оптимальность методов оптимизации	32
2.4. Теоретические основы сходимости одномерных алгоритмов глобального поиска	35
2.5. Индексная схема учета ограничений	45
Краткое содержание главы	47
Глава 3. Фундаментальные способы редукции размерности. Многошаговая схема	49
3.1. Принципы редуцирования сложности в задачах принятия решений.	49
3.2. Многошаговая схема редукции размерности	51
3.3. Свойства одномерных подзадач многошаговой схемы	59
3.3.1. Структура допустимых областей одномерного поиска	59
3.3.2. Свойства целевых функций в одномерных подзадачах	64
Краткое содержание главы	65
Глава 4. Модели и методы поиска локально-оптимальных решений	65
4.1. Постановка задачи поиска локально-оптимальных решений	65
4.2. Общие принципы построения методов локальной оптимизации	66
4.2.1. Структура методов поиска локального минимума функций	66
4.2.2. Измерения локальной информации и роль модели задачи в их интерпретации	67
4.2.3. Классификация методов локального поиска	68
4.2.4. Эффективные стратегии поиска вдоль направлений. Регуляризованные алгоритмы одномерного поиска	69
4.3. Классические методы локальной оптимизации	71
4.4. Методы локальной оптимизации, основанные на квадратичной модели поведения функций	76
4.4.1. Методы второго порядка для гладких задач	77
4.4.2. Методы первого порядка для гладких задач	83
4.5. Некоторые методы прямого поиска для негладких задач	94
4.5.1. Метод Нелдера–Мида	95

4.5.2. Метод Хука-Дживса	96
4.6. Особенности применения методов локального поиска при двусторонних ограничениях на переменные	98
4.6.1. Особенности учета двусторонних ограничений на переменные в методах гладкой оптимизации	98
4.6.2. Учет двусторонних ограничений в методах прямого поиска	101
4.7. Учет ограничений общего вида на основе метода штрафов	101
4.7.1. Метод внешнего штрафа. Общие условия сходимости	102
4.7.2. Структура возникающих задач со штрафом и характер приближения оценок к решению	106
4.7.3. Недостаточность локальных методов при использовании метода штрафов	109
4.7.4. Сочетание локальных методов с методами покрытий области	110
Краткий обзор главы	110
Контрольные вопросы и упражнения	111
Предметный указатель	111
Литература	115

Глава 1. Математические модели принятия решений. Примеры прикладных задач

1.1. Методология принятия решений

Задачи минимизации (или *максимизации*) функций при различных дополнительных условиях являются типичными математическими моделями процессов выбора решений при автоматизированном проектировании технических устройств и систем, в управлении подвижными частями роботов, при восстановлении зависимостей на основе анализа экспериментальных данных и т.д. [1,2]

Простейший тип такой задачи предполагает, что выбор варианта объекта оптимизации характеризуется выбором значения *вектора варьируемых параметров* $y = (y_1, y_2, \dots, y_N)$. При этом условия создания и функционирования объекта накладывают некоторые ограничения на допустимые значения вектора y , которые формально описываются как требования принадлежности вектора y некоторой *допустимой области* Q в N -мерном пространстве параметров R^N . Эффективность варианта объекта, соответствующего заданному значению вектора y , описывается *показателем* $f(y)$, который может быть вычислен на основе анализа математической модели этого объекта. Такие вычисления будем называть *испытаниями*. Если принять, что уменьшение показателя $f(y)$ соответствует улучшению проекта, то выбор значения y^* , определяющего наилучший вариант, сводится к приближенному решению задачи минимизации

$$y^* = \arg \min \{f(y) : y \in Q\} \quad (1.1)$$

на основе результатов $Z^i = f(y^i), 1 \leq i \leq k$ конечного числа испытаний в точках $y^i \in Q$.

Выбор точек испытаний может осуществляться последовательно, т.е. при выборе очередной точки y^{i+1} могут быть использованы уже известные результаты испытаний Z^1, \dots, Z^i в предшествующих точках y^1, \dots, y^i . Допускается, что некоторые выбранные точки могут не принадлежать области Q . Математическая модель объекта должна содержать тест для выявления таких случаев (при этом отрицательный результат теста принимается за исход испытания). Таким образом, процедура выбора точек из области Q может быть определена косвенно – через тест на принадлежность к этой области. Возможен и случай, когда допустимая область пуста, что соответствует несовместимости требований, предъявляемых к объекту. Установление этого факта также считается решением задачи (1.1).

Точность приближенного решения задачи (1.1), имеющего вектор варьируемых параметров y^* и значение $z^* = f(y^*)$, может характеризоваться некоторым принятым (из содержательных соображений) понятием близости к точному решению.

Многоэкстремальность функции f приводит к тому, что вычислительные схемы, основанные на наглядных соображениях типа геометрической идеи “скатывания по склону” до ближайшего минимального значения, могут приводить в локальный минимум, не гарантируя получения оценок точки y^* абсолютного минимального значения. Идеи, на которых основаны решающие правила алгоритмов, гарантирующих оценки глобального минимума, достаточно абстрактны [20, 25, 26]. Типичная схема состоит в том, что минимизируемая функция рассматривается как представитель некоторого класса функций, зависящего от параметра. Например, предполагается, что f есть липшицева функция (с константой L) или что f есть реализация винеровского случайного процесса и т.п. При этих предположениях предлагается или выводится

решающее правило алгоритма и обосновывается его применимость для получения оценок величин y^* и $z^*=f(y^*)$.

Правила, соответствующие конкретному алгоритму, порождают для конкретной задачи последовательность точек $\{y^k\}$, в которых выполняются итерации (называемую также минимизирующей последовательностью). Условие остановки делает эту последовательность конечной, причем любому ее усечению $\{y^0, \dots, y^k\}$ соответствует усеченная последовательность вычисленных значений функции, по которой строится текущая оценка искомого минимального значения.

Теоретический вывод (если он имеет место) или эвристическая мотивация соответствующего решающего правила прямо или косвенно направлены на то, чтобы алгоритм обеспечивал достижение заданной точности оценок оптимума при возможно меньшем числе итераций. С этой целью очередная итерация производится, например, в точке минимума нижней огибающей класса липшицевых функций, построенной с учетом результатов выполненных итераций, или в точке минимума условного математического ожидания значения винеровского процесса и т.п.

Точки итераций, порождаемые правилами подобного рода, характеризуются неравномерным расположением в области поиска, сгущаясь в окрестностях глобального и близких к нему (по значению минимизируемой функции) локальных минимумов. Чем более неравномерно расположение этих точек (при заданном числе итераций) и чем меньше обеспечиваемый за счет этой неравномерности радиус окрестности точки глобального минимума, тем эффективнее метод. Поэтому теоретические свойства таких алгоритмов обычно описываются в форме сравнения оценок плотности итераций в различных частях интервала поиска с плотностью в окрестности глобального минимума. При асимптотическом рассмотрении, т.е. в случае неограниченного продолжения поиска, соответствующие оценки устанавливают связь предельных точек последовательности с точками глобального минимума.

1.2. Общая модель рационального выбора

Рассматривается общая задача оптимизации, охватывающая основные постановки задач выбора оптимального решения [26].

Математическая модель объекта оптимизации содержит вектор варьируемых параметров

$$y = (y_1, y_2, \dots, y_N) \quad (1.2)$$

вектор-функцию характеристик

$$W(y) = (w_1(y), \dots, w_n(y)), \quad (1.3)$$

причем все координатные функции $w_i(y)$ таковы, что их уменьшение соответствует улучшению проекта.

Значение вектора параметров y , характеризующее проектное решение, должно принадлежать гиперинтервалу (области поиска)

$$D = \{y \in R^N : a_i \leq y_i \leq b_i, 1 \leq i \leq N\}, \quad (1.4)$$

определяемому заданными векторами начала и конца

$$a = (a_1, \dots, a_N), \quad b = (b_1, \dots, b_N) \quad (1.5)$$

где

$$a_i < b_i, \quad 1 \leq i \leq N. \quad (1.6)$$

Часть координатных функций $w_i(\mathbf{y})$ из (1.3) выбираются в качестве критериев эффективности требуемого решения. С этой целью задается множество

$$F = \{i_1, \dots, i_s\} \subset \{1, \dots, n\}, \quad (1.7)$$

включающее номера таких координатных функций. В результате определяется *векторный критерий эффективности*

$$\mathbf{f}(\mathbf{y}) = (f_1(\mathbf{y}), \dots, f_s(\mathbf{y})) = W_F(\mathbf{y}) \quad (1.8)$$

где

$$f_j(\mathbf{y}) = w_{i_j}(\mathbf{y}), i_j \in F. \quad (1.9)$$

При выборе проектного решения желательно обеспечить возможно большее уменьшение каждой из координатных функций, входящих в векторный критерий (1.8) и называемых частными критериями.

Номера координатных функций из (1.3), не вошедших в векторный критерий эффективности (1.8), образуют множество

$$G = \{j_1, \dots, j_m\} = \{1, \dots, n\} \setminus F, \quad (1.10)$$

где

$$m = n - s. \quad (1.11)$$

Функции $w_j(\mathbf{y}), j \in G$ требуется минимизировать до выполнения неравенств

$$w_{j_i}(\mathbf{y}) \leq q_i, j_i \in G \quad (1.12)$$

где все величины $q_i, 1 \leq i \leq m$ из правых частей неравенств являются заданными и образуют *вектор допусков*

$$\mathbf{q} = (q_1, \dots, q_m), \quad (1.13)$$

причем

$$q_i > 0, \quad 1 \leq i \leq m. \quad (1.14)$$

Неравенства (1.12) выделяют множество (*область допустимых решений*)

$$Q = \{\mathbf{y} \in D : w_{j_i}(\mathbf{y}) \leq q_i, \quad 1 \leq i \leq m\}, \quad (1.15)$$

содержащее допустимые проектные решения \mathbf{x} .

Обозначения

$$g_i(\mathbf{y}) = w_{j_i}(\mathbf{y}) - q_i, \quad 1 \leq i \leq m \quad (1.16)$$

позволяют преобразовывать условия (1.12) в эквивалентную форму

$$g_i(\mathbf{y}) \leq 0, \quad 1 \leq i \leq m, \quad (1.17)$$

в дальнейшем называемую *функциональными ограничениями* задачи оптимизации.

Левые части ограничений (1.17) образуют вектор ограничений

$$\mathbf{g}(\mathbf{y}) = (g_1(\mathbf{y}), \dots, g_m(\mathbf{y})) = W_G(\mathbf{y}) - \mathbf{q}, \quad (1.18)$$

неположительность координатных функций которого является необходимым и достаточным условием допустимости решения \mathbf{y} .

Схематично математическая модель объекта оптимизации представлена на рис.1.1.

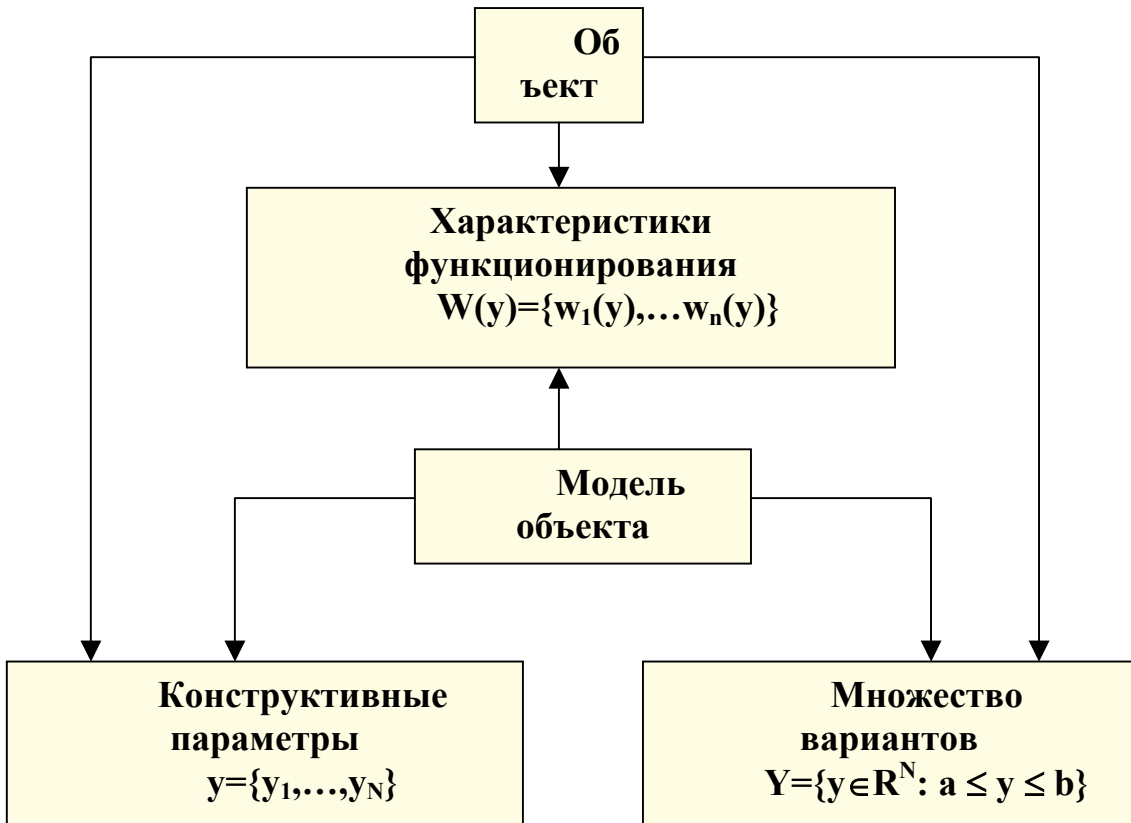


Рис. 1.1. Математическая модель объекта оптимизации

Требования к рациональному варианту проектируемого объекта представлены на рис.1.2.

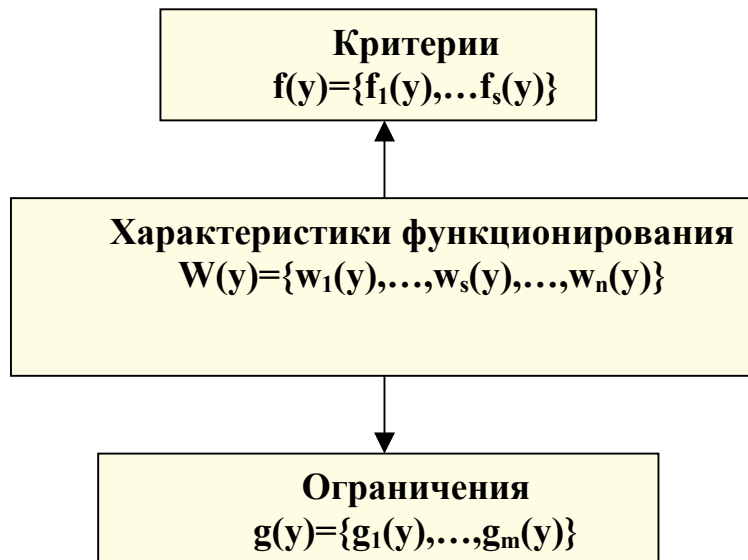


Рис. 1.2. Рациональность варианта проектируемого объекта

1.3. Модели математического программирования и методы их анализа

Тип задач принятия оптимального решения определяется следующими основными свойствами задач [38]:

- число критериев эффективности (однокритериальная (скалярная) / многокритериальная задача);
- наличие ограничений;
- число экстремумов (одноэкстремальная / многоэкстремальная задача);
- свойства функций критериев эффективности (линейность, нелинейность, непрерывность, разрывность и т.д.);
- свойства функций ограничений;
- свойства управляемых параметров (непрерывность, дискретность);
- свойства области допустимых решений (односвязность, многосвязность).

В области вычислительных методов поддержки принятия оптимальных решений можно выделить группы задач, где в настоящее время уже получены важные теоретические результаты, разработаны методы, изданы монографии и учебники. В то же время остаются задачи, исследования которых не завершены. Часть результатов по ним представлена в оригинальных монографиях, а часть доступна только через публикации в научных журналах.

В следующем кратком обзоре научных направлений будут указаны как классические, так и развивающиеся в настоящее время направления, но большее внимание будет уделено последним. В обзор не включены направления, связанные со специальными классами задач: линейными, квадратичными, выпуклыми.

Один из наиболее завершенных разделов образует класс *методов локальной оптимизации без ограничений*. Сюда относятся классические результаты по методам прямого поиска Нелдера и Мида (1965 г.), Хука и Дживса (1961 г.), результаты по обширной группе методов, основанных на квадратичных моделях поведения оптимизируемых функций. Здесь можно назвать работы периода 1965-1972 гг. Бройдена С., Флетчера Р., Пауэлла М., Шанно Д., Гольдфарба Д., Давидона С., Ривса С., Данилина Ю.М., Поляка Б.Т., Пшеничного Б.Н., Шора Н.З. и других авторов. Результаты отражены во многих монографиях [5, 7, 8, 12, 14] и учебниках [6, 9, 11, 13, 19].

Результаты этого раздела широко опираются на методы математического анализа и линейной алгебры. Разработанные методы обобщаются на задачи с простыми линейными ограничениями.

Следующее направление, породившее в 60-х, 70-х годах значительное количество работ, связано с *общими и специальными методами учета нелинейных ограничений* в задачах оптимального выбора. Общие методы основываются на использовании штрафных и барьерных функций, модифицированных функций Лагранжа, а также метода центров и параметризации целевой функции (последняя группа методов приведена в [22]). При этих подходах задача со сложными ограничениями трансформируется в последовательность вспомогательных задач без таких ограничений с применением к ним известных методов безусловной оптимизации. Однако, возникающие вспомогательные задачи обычно имеют плохую структуру. Поэтому большее количество исследований было посвящено попыткам построения специальных методов локальной оптимизации с учетом ограничений. Результаты всех этих исследований достаточно полно отражены в литературе. Укажем на монографии [7, 10, 12, 14-18, 21] и учебники [6, 9, 15, 13, 19].

Значительная часть упомянутых методов локального поиска формирует направления локального убывания функции и выполняет смещение вдоль них, используя *специальные процедуры одномерной локальной оптимизации*. Многократное повторение этих процедур требует их оптимизации. Здесь широко используются результаты теории построения оптимальных алгоритмов [3, 4]. В учебной литературе эти вопросы освещены в [6-9]. Важное место занимают вопросы регуляризации решения некорректных задач и *метод регуляризации*, предложенный Тихоновым А.Н. [19].

Особое место занимают научные *исследования в области многоэкстремальной оптимизации*. Поскольку глобальный экстремум является интегральной характеристикой функций задачи в области поиска, его определение связано, в общем случае, с построением покрытия области точками испытаний.

Применяемые подходы существенно зависят от имеющейся априорной информации о задаче и затрат, связанных с однократным вычислением функций задачи. Все методы приводят к явному или неявному построению покрытий области. Простейшие методы решения связаны с применением *неадаптивных случайных или детерминированных равномерных покрытий*. Интересные решения задач о построении равномерных регулярных многомерных решеток связаны с работами Соболя И.М. [22], пионерские работы по применению стохастического подхода принадлежат Растригину Л.А. [21]. Разнообразные результаты по глобальному случайному поиску представлены в монографии Жиглявского А.А. [23]. Оригинальный подход на основе стохастического оценивания меры областей со значениями минимизируемой функции меньшими некоторых пороговых значений представлен в монографии Чичинадзе В.К. [24]. Все эти подходы почти не используют априорную информацию о решаемой задаче и обычно требуют большого объема вычислений для достоверного оценивания решения.

Альтернативное направление научных исследований основано на *оптимальном планировании размещения точек испытаний* с использованием имеющейся априорной, а также получаемой в ходе поиска информации о решаемой задаче. В основе этого направления лежит понятие модели решаемой задачи. Метод решения строится как оптимальное, в том или ином смысле, решающее правило в рамках подходов теории игр и оптимальных статистических решений. Исследованием вопросов построения подобных вычислительных процедур занимались: Кушнер Х., Неймарк Ю.И., Стронгин Р.Г., Сухарев А.Г., Евтушенко Ю.Г., Моцкус Й.Б., Шалтянис В.Р., Жилинскас А.Г., Пиявский С.А., Хансен П., Пинтер Я., Пардалос П.М., Хорст Р. и многие другие исследователи, входящие в состав связанных с этими работами научных школ. Методы, разрабатываемые в рамках этого направления, порождают *неравномерные адаптивные покрытия области поиска* точками измерений, существенно более плотные в окрестностях решения. Большинство методов, в рамках данного подхода, строится как одношагово-оптимальные правила. При этом используются или гарантирующие модели поведения, основанные на построении минорант функций (Евтушенко Ю.Г., Пиявский С.А., Хансен П., Пинтер Я.) или вероятностные модели поведения, когда функция рассматривается как реализация некоторого случайного процесса (Кушнер Х., Неймарк Ю.И., Стронгин Р.Г., Моцкус Й.Б., Шалтянис В.Р., Жилинскас А.Г.) или аксиоматические модели функции (Жилинскас А.Г.).

Наиболее полные результаты получены в одномерных задачах многоэкстремальной оптимизации. Для указанного класса методов глобальной оптимизации получены *новые результаты для задач с ограничениями*. В рамках информационно-статистического метода, предложенного Р.Г.Стронгиным, разработан *высокоэффективный индексный алгоритм учета ограничений* (результат получен Маркиным Д.Л.) [26]. Индексный метод обладает рядом преимуществ по сравнению с

традиционными методами: не использует штрафных добавок, не требует подбора коэффициентов типа констант штрафа и решения последовательностей безусловных подзадач, допускает частичную вычислимость целевой функции и ограничений.

При разработке алгоритмов важный теоретический фундамент образуют работы по аналитическому исследованию сходимости итерационных последовательностей к локальным и глобальным решениям многоэкстремальных задач. Возможен общий взгляд на широкие классы методов глобальной оптимизации, позволяющий исследовать условия и характер их сходимости в весьма общей форме. Он основан на модели *характеристической представимости алгоритмов*, впервые предложенной Гришагиным В.А. (ее описание можно найти, например, в [35]), и ее обобщениях, применимых для анализа сходимости специальных многомерных вычислительных схем.

Более сложным, по сравнению со случаем одной переменной, является построение одношагово-оптимальных процедур глобальной оптимизации для многомерных задач. В настоящее время развивается несколько направлений. Первое основано на *редукции размерности задачи*. При одном из подходов многомерная задача сводится к серии вложенных одномерных задач (*многошаговая схема* редукции размерности). Остальные подходы, разработанные в рамках информационно-статистических алгоритмов, используют *отображение многомерной области на отрезок с помощью разверток* на основе аппроксимаций кривых Пеано [25, 26]. Последние результаты в этой области получены Р.Г.Стронгиным и связаны с *применением множественных разверток* [28], позволяющих более точно передать близость точек измерений в многомерном пространстве при его отображении на отрезок.

Другое направление основано на *адаптивном разбиении области поиска на подобласти - компоненты простой структуры*, обычно – параллелепипеды [29-31]. Для каждой подобласти учитываются только выполненные в ней измерения функций задачи, что упрощает получение оценок поведения функции в каждой такой компоненте. На основе этих оценок компоненте приписывается приоритет. На каждой итерации наиболее приоритетная компонента разделяется на несколько новых компонент, с проведением в них дополнительных вычислений функции. В настоящее время в рамках этого направления разрабатываются новые методы.

Еще одним направлением исследований является *учет различной информации, получаемой о функции в результате измерения*, например, градиента, а также распространение разработанных методов на новые классы функций. Имеются относительно новые результаты, связанные с решением сложных *многокритериальных задач*. Наличие нескольких критериев приводит к изменению понятия решения, которое теперь понимается в смысле решения по Парето [32, 33]. Известные методы свертки [2] не позволяют одновременно оценивать это множество в целом. Один из последних результатов в этой области состоит в том, что удалось построить задачу со скалярным перестраиваемым критерием, минимизация которого приводит к *оцениванию сразу всего множества Парето*. Этот результат получен в работе [36].

Интенсивные исследования проводились в области *параллельной глобальной оптимизации* применительно к классу информационно-статистических алгоритмов. Результаты представлены в [28,35] .

Учитывая новизну и интенсивность проводимых научных исследований, центральную часть настоящего учебного пособия будут составлять эффективные методы, разработанные для поиска решения многоэкстремальных задач с невыпуклыми ограничениями. Многие оригинальные результаты в этой области принадлежат нижегородской школе глобальной оптимизации.

1.4. Принципы создания учебно-исследовательских систем принятия решений

Системы имитации объектов и явлений на компьютере основаны на использовании математических моделей этих объектов и явлений. Поэтому образовательное использование таких систем ограничивается степенью адекватности, обеспечиваемой соответствующей математической моделью. В пределах адекватности, гарантируемой моделью, имитационная система дает новые возможности восприятия и рождает новые стимулы к познанию по сравнению с традиционными подходами. Имитация на компьютере позволяет наблюдать динамику объекта изучения в темпе, характерном для человеческого восприятия, хотя подлинные времена течения процессов могут составлять доли секунды (взрыв) или годы (движение горных пород). Машинные средства визуализации позволяют создавать наглядные образы объектов и явлений, которые сами по себе не являются наглядными. При этом создаваемые образы правильно отражают основные отношения, характеризующие исходный объект. Средства имитационной системы могут адаптироваться к скорости реакции, наблюдательности, другим возможностям конкретного обучаемого. Разыгрывание на компьютере различных вариантов и сравнение результатов выбора создают поле для самостоятельных выводов и развития интуиции. Возможность придания этому процессу индукции соревновательного характера является дополнительным познавательным стимулом. Имитатор защищает обучаемого от раздражающих и отвлекающих технических ошибок (неверное нажатие клавиш, неточное задание условий и т.п.), позволяя сосредоточиться на изучаемой теме. Эти черты имитационной системы создают образовательную среду, интенсифицирующую индуктивную и дедуктивную активность обучаемого.

Проектирование имитационных систем, обладающих рассмотренными образовательными возможностями, начинается с разработки *основных компонент сценария обучения* и исследований в среде такой системы [37]. При этом исследовательский аспект имеет первостепенное значение, ибо исследование есть одновременно одна из важнейших форм и одна из важнейших целей обучения. Основные вопросы, возникающие при определении главных компонент и целей сценария, коротко состоят в следующем (рис. 1.3).

Существует широкий круг тем, изучение которых может ускоряться и углубляться благодаря высокой степени наглядности, обеспечиваемой имитатором. При этом *объектом изучения* может быть явление или понятие, метод исследования или расчетов, конкретный прибор и т.п. Каждая такая тема предполагает существование или разработку соответствующей математической модели, отражающей изучаемые аспекты темы и их взаимосвязи. Центральный вопрос, возникающий после выбора изучаемой темы и математической модели, поддерживающей это изучение, состоит в выборе *объекта показа*, многовариантная имитация которого раскрывает тему. Объект показа, будучи целым и единым, состоит из взаимосвязных частей и сторон, характеристики которых должны быть измерены и продемонстрированы, для чего нужно определить датчики свойств объекта и наглядные образы для визуализации этих свойств на дисплее.

Освоение темы, т.е. осознание законов природы, теорем, гипотез, принципов функционирования устройств и т.п. опирается на анализ многих вариантов, выбор и сопоставление которых в режиме диалога обучаемого и компьютера должен предусматривать сценарий. Сценарий должен обеспечивать раскрытие взаимосвязи прикладных и фундаментальных аспектов темы, ибо, с одной стороны, понимание любого частного положения опирается на фундаментальные представления, а с другой стороны, освоение фундаментальных представлений достигается и углубляется при

знакомстве с их частными приложениями. Поэтому образовательная имитационная система по любой частной теме должна служить и задачам усвоения и углубления фундаментальных представлений.

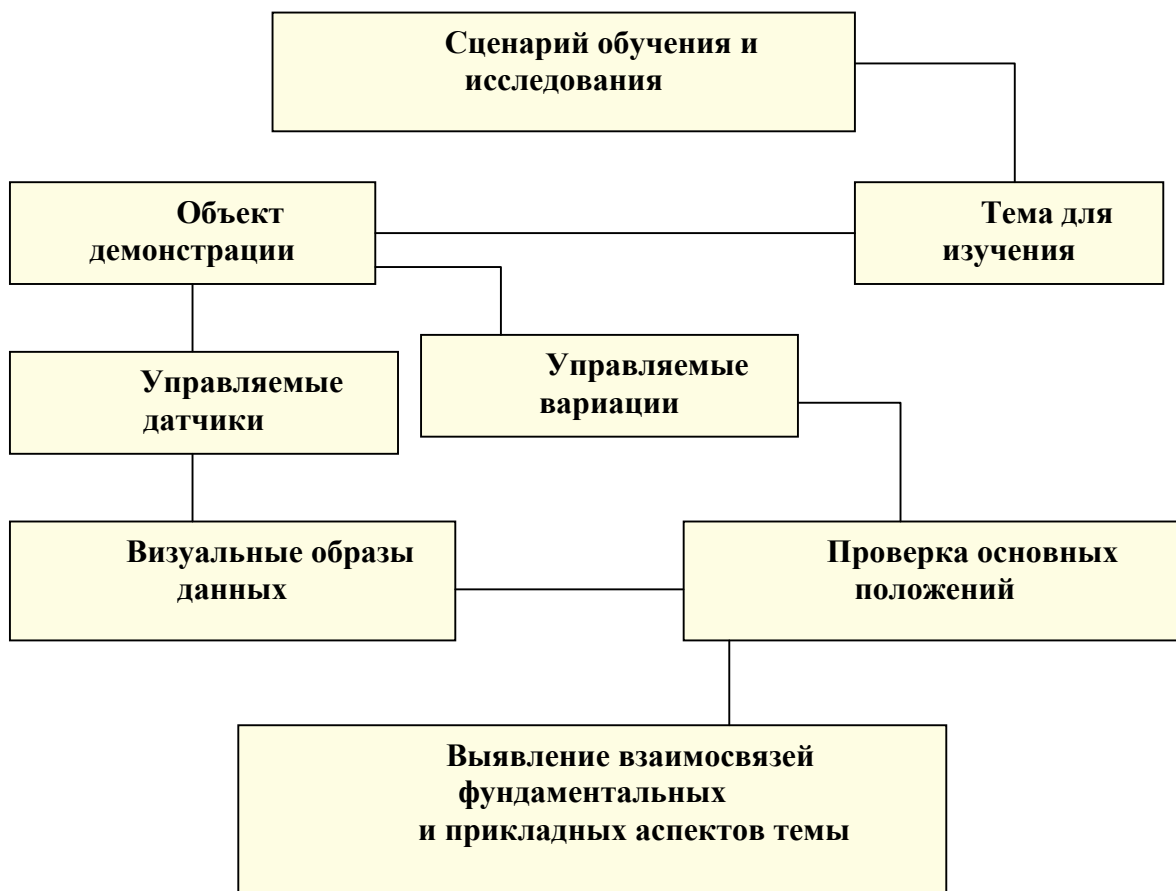


Рис. 1.3. Основные компоненты и цели сценария

В качестве изучаемой темы, рассмотрим численные методы для поиска абсолютного минимума многоэкстремальной функции. Свойства минимизирующих последовательностей отражают как характер исходных предположений о классе минимизируемых функций (включая зависимость от параметра, входящего в описание класса: константа Липшица, параметр винеровского процесса и т.п.), так и принципы построения решающего правила, условия останова и прикладные возможности алгоритма. Поэтому демонстрация минимизирующих последовательностей, порождаемых конкретными правилами для конкретных задач, наглядно проявляющая обсуждаемые связи, ускоряет и углубляет восприятие теории глобальной оптимизации, развивает интуицию, необходимую для практического использования (и дальнейшего развития) соответствующих методов.



Рис. 1.4. Основные компоненты и цели сценария в системе АБСОЛЮТ

Графический образ усечения минимизирующей последовательности $\{y^k\}$ создается с помощью набора вертикальных штрихов, отмечающих координаты точек y^k , в которых проводились вычисления значений функции (см. линейчатый "спектр", изображенный на рис.1.5; данный рисунок получен при помощи учебно-исследовательской системы АБСОЛЮТ[37]). Этот набор дополняется ломаной линией, характеризующей порядок выполнения итераций. Отрезки указанной ломаной линии соединяют пары точек вида (k, y^k) , $(k+1, y^{k+1})$, где y^k есть координата точки, в которой на k -той итерации вычислялось значение $z^k=f(y^k)$ (см. изображение ломаной на рис.1.5 под линейчатым спектром). Аналогично создается графический образ последовательности $\{z^k\}$.

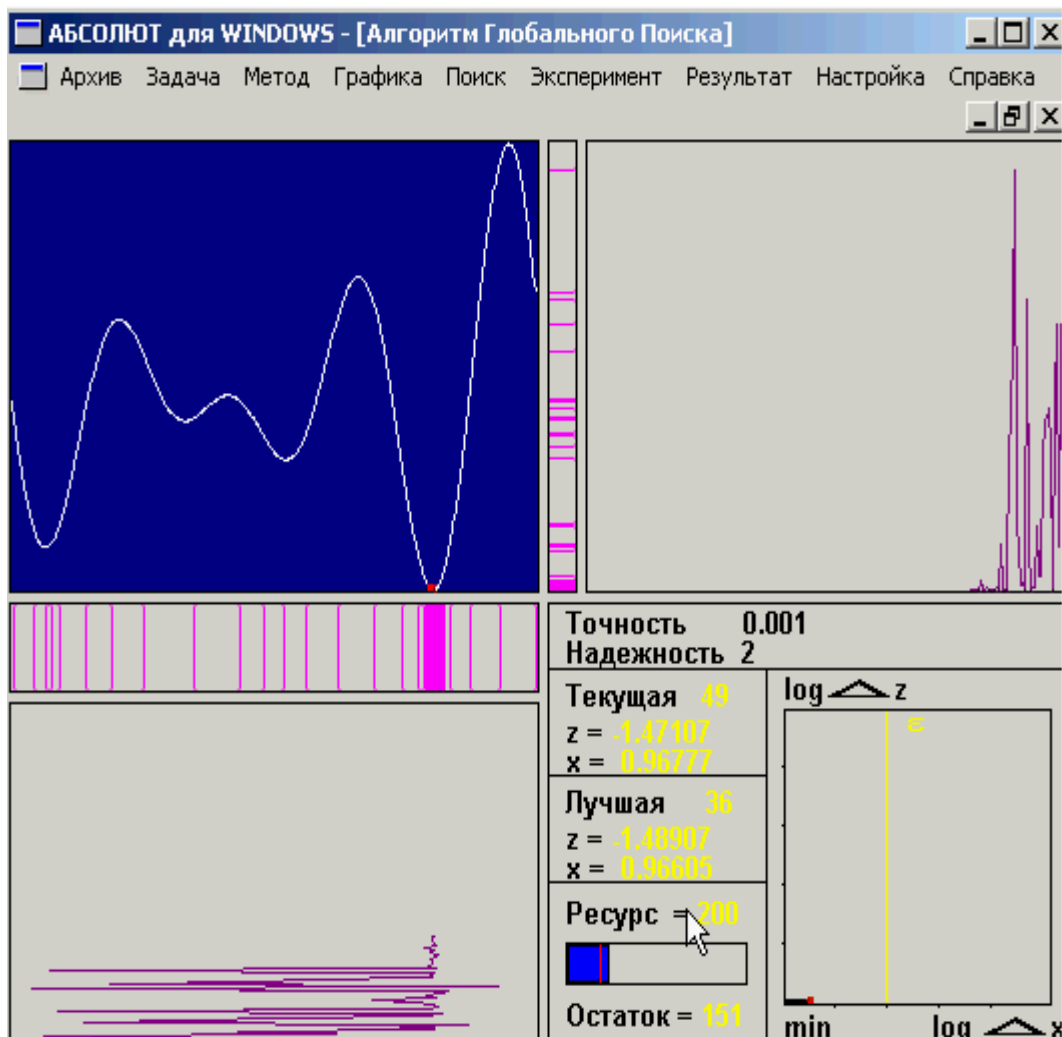


Рис. 1.5. Графическое представление одномерной оптимизации

Возможность анализа локального строения последовательности (например, в окрестностях точек накопления, где плотность итераций резко возрастает) обеспечивается “увеличительным стеклом”, растягивающим выделенную окрестность в заданное число раз. Настройку стекла, т.е. указание окрестности и масштаба изображения, пользователь осуществляет самостоятельно. Этот механизм настраиваемой лупы используется и при ручном выборе точек итераций.

Общую характеристику распределения итераций в интервале $[a,b]$ дают гистограммы, автоматически перевычисляемые на каждом шаге поиска минимума. При этом ломаная линия, характеризующая последовательность итераций, высвечивается на фоне соответствующей гистограммы (с использованием различных цветов изображений). Наглядное сравнение последовательностей, сопоставляемых различными алгоритмами (или одним алгоритмом, но при различных значениях параметров), обеспечивается одновременной демонстрацией нескольких изображений.

Для изучения эффективности условий останова, входящих в описание алгоритмов, система снабжена логарифмическим индикатором точности, демонстрирующим логарифм отклонения текущей оценки оптимума от истинного значения. Это позволяет наглядно наблюдать как случаи, когда остановка поиска происходит по достижению заданной точности, так и случаи продолжения итераций в ситуации, когда фактически достигнутая точность превышает заданную.

Поскольку характер последовательностей, порождаемых конкретными алгоритмами, зависит от свойств минимизируемых функций, то для демонстрации различных вариантов такой зависимости система должна предоставлять тестовые функции, либо выбираемые из стандартного набора, либо задаваемые формульным описанием, либо генерируемые случайным механизмом.

1.5. Прикладные задачи оптимального выбора

1.5.1. Задача определения местоположения и параметров движения объекта по измерениям высоты рельефа местности

Данная задача относится к группе задач навигации по геофизическим полям, рассмотренным, например в [45].

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 Q его возможных значений

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

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

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

где ζ_i – независимые реализации центрированной составляющей помехи измерений с плотностью распределения $P_\zeta(z)$, 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_n и результатам h_1, h_2, \dots, h_n измерений высот рельефа вдоль траектории полета оценить координаты y_1, y_2 начального местоположения объекта в области Q .

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 наблюденных значений высот рельефа при условии, что начальное местоположение и систематическая составляющая измерений имеют значения y_1, y_2, C .

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

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

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

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

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

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

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

где

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

При этом минимум по C можно найти аналитически из условия $\frac{\partial f}{\partial C}(y_1, y_2, C) = 0$.

Отсюда находим

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

Окончательно, задача определения местоположения сводится к следующей задаче многоэкстремальной оптимизации

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

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

Многоэкстремальный характер функции (1.26) связан с существованием на карте местности участков с похожими сечениями рельефа. Постановка задачи в виде (1.25), (1.26) соответствует случаю нормальной помехи измерений. В действительности эта гипотеза верна не всегда. При измерениях возможны редкие большие случайные отклонения не соответствующие нормальному закону. Более адекватной моделью помехи является предположение, что распределение $P_\xi(z)$ неизвестно, но принадлежит известному классу распределений \mathbf{P} . В этом случае применяют огрубленный метод максимального правдоподобия, состоящий в следующем [46].

В классе \mathbf{P} нужно выбрать распределение $P_\xi^*(z)$, на котором достигает максимума функционал информации Фишера

$$I(p) = \int (P_\xi'(z))^2 / P_\xi(z) dz, \quad (1.27)$$

а затем в качестве распределения $P_\xi(z)$ в (1.21), (1.22) использовать найденное наименее благоприятное распределение $P_\xi^*(z)$. Например, если \mathbf{P} – класс произвольных “невырожденных” распределений, для которых $P_\xi(0)$ равномерно отделено от нуля $P_\xi(0) \geq 1/2d > 0$, то наименее благоприятным является распределение Лапласа

$$P_\xi^*(z) = 1/2d e^{-\frac{|z|}{d}} .$$

Следовательно, в этом классе помех в задаче оценивания (1.23) следует вместо функции (1.24) использовать

$$F(y_1, y_2, C) = 1/n \sum |h_i - h(y_1(t_i), y_2(t_i)) - C|. \quad (1.28)$$

с учетом оптимального выбора значения C приходим к задаче оценивания двух переменных y_1 и y_2 следующего вида

$$f(y_1, y_2) \rightarrow \min, (y_1, y_2) \in Q, \tag{1.29}$$

$$f(y_1, y_2) = \min_{j=1, \dots, n} \frac{1}{n} \sum_{i=1}^n \left| (h_i - h(y_1 + v_1 t_i, y_2 + v_2 t_i)) - (h_j - h(y_1 + v_1 t_j, y_2 + v_2 t_j)) \right|.$$

1.5.2. Задача определения обобщенных координат манипулятора для заданного позиционирования схвата

Пусть имеется манипулятор, включающий N+1 звено. Звеньям присвоим номера 0, 1, 2, ..., N. Звено с номером 0 неподвижно. Остальные звенья последовательно связаны друг с другом и образуют кинематическую цепочку, в которой каждое следующее звено имеет ровно одну относительную степень подвижности по отношению к предыдущему звену. В этом случае принято говорить, что каждые два последовательно взятые в кинематической цепочке звена образуют кинематическую пару пятого класса. Класс кинематической пары определяется как разность числа обобщенных координат у свободного твердого тела, равного шести, и числа относительных степеней подвижности в кинематической паре. Обобщенные координаты, описывающие относительное смещение звеньев в кинематических парах образуют вектор $q = (q_1, q_2, \dots, q_N)$. Будем считать, что все звенья являются абсолютно жесткими твердыми телами с известными геометрическими характеристиками. Каждому набору обобщенных координат q_1, q_2, \dots, q_N будет однозначно соответствовать пространственная конфигурация манипулятора, а также конкретное положение и ориентация в пространстве его последнего звена – схвата.

Содержательно задача состоит в том, чтобы при имеющейся информации о виде манипулятора определить те значения его обобщенных координат, при которых схват манипулятора получает конкретное требуемое положение и ориентацию в пространстве.

Для того чтобы поставить эту задачу математически нужно ввести формальное описание кинематики манипулятора. Введем для этого ряд соглашений.

Будем рассматривать такие манипуляторы, у которых относительное перемещение звеньев в кинематических парах состоит либо в поступательном смещении вдоль некоторого направления, либо в повороте вокруг некоторой оси. Такие пары принято называть кинематическими парами поступательного или вращательного типов. Ось O_i , вдоль которой происходит смещение или поворот называют осью i -той кинематической пары.

На рис. 1.6 изображены примеры двух пар вращательного и одной пары поступательного типа.

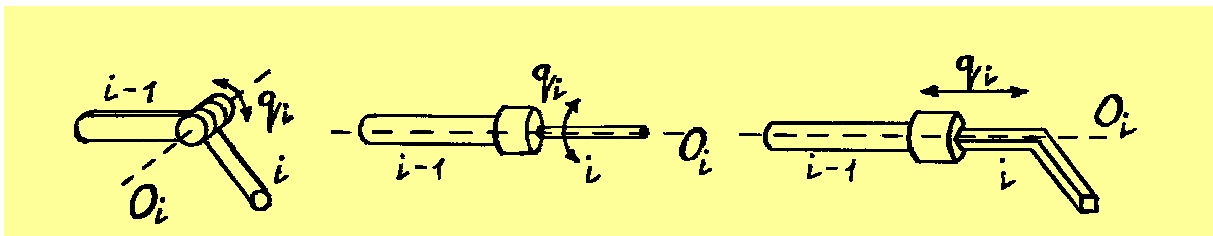


Рис. 1.6. Примеры кинематических пар

Для описания положения и ориентации звеньев манипулятора в пространстве введем системы координат, жестко связанные с этими звеньями.

С каждым i -м звеном свяжем две правые системы координат, одну – с его началом, другую – с его концом. Эти системы будем называть локальными. Первую

обозначим $e_{x_i}, e_{y_i}, e_{z_i}$ и будем называть системой (i) , а вторую обозначим $\hat{e}_{x_i}, \hat{e}_{y_i}, \hat{e}_{z_i}$ и будем называть системой (\hat{i}) . Базисные векторы e_i будем считать ортонормированными. Центры этих систем размещаются на осях кинематических пар.

Система координат $e_{x_0}, e_{y_0}, e_{z_0}$, связанная с началом неподвижного нулевого звена, в дальнейшем будет рассматриваться как абсолютная система координат, по отношению к которой будет определяться положение и ориентация последнего звена – схвата.

Взаиморасположение систем координат подчиним определенным правилам. Для этого введем некоторую стандартную начальную пространственную конфигурацию манипулятора, условно соответствующую нулевым значениям обобщенных координат.

Установим для этой начальной конфигурации следующие правила взаиморасположения локальных систем координат. Их размещение начнем с произвольного выбора в начале нулевого неподвижного звена системы (0) . Далее размещается система $(\hat{0}), (1), (\hat{1})$ и т.д.

Пусть в начале i -го звена уже размещена локальная система координат $e_{x_i}, e_{y_i}, e_{z_i}$. Тогда система (\hat{i}) размещается в конце i -го звена так, что ось \hat{e}_{z_i} выбирается в направлении оси O_{i+1} , а ось e_{x_i} ортогональна к осям e_{z_i} и \hat{e}_{z_i} . Если эти оси параллельны, то указанное правило не позволяет однозначно построить ось \hat{e}_{x_i} . Для устранения неоднозначности, \hat{e}_{x_i} выбирается параллельной оси e_{x_i} . Ось \hat{e}_{y_i} всегда дополняет систему \hat{i} до правой системы координат.

После того, как система (\hat{i}) , жестко связанная с концом i -го звена, построена, с началом $i+1$ звена связывается система $(i+1)$, совпадающая при стандартной конфигурации манипулятора по своему положению и ориентации с системой (\hat{i}) . Отличие этих систем в том, что они связаны с разными звеньями. А именно, система $(i+1)$ жестко связана с $(i+1)$ звеном. Поэтому, при изменении конфигурации манипулятора в пространстве система $(i+1)$ начнет изменять свое положение относительно системы (\hat{i}) . Изменение положения будет происходить либо за счет ее смещения вдоль оси \hat{e}_{z_i} на величину q_i (для кинематических пар поступательного типа), либо за счет поворота вокруг \hat{e}_{z_i} на угол q_i (для кинематических пар вращательного типа).

Последнюю систему координат \hat{N} , связанную с концом схвата стандартно выбирают следующим образом. Центр ее помещают между губками схвата. Ось \hat{e}_{z_N} выбирают вдоль линии действия схвата, ось \hat{e}_{x_N} «протыкают» губки схвата, а ось \hat{e}_{y_N} дополняют систему до правой. Выбор этой системы координат показан на рис. 1.7.

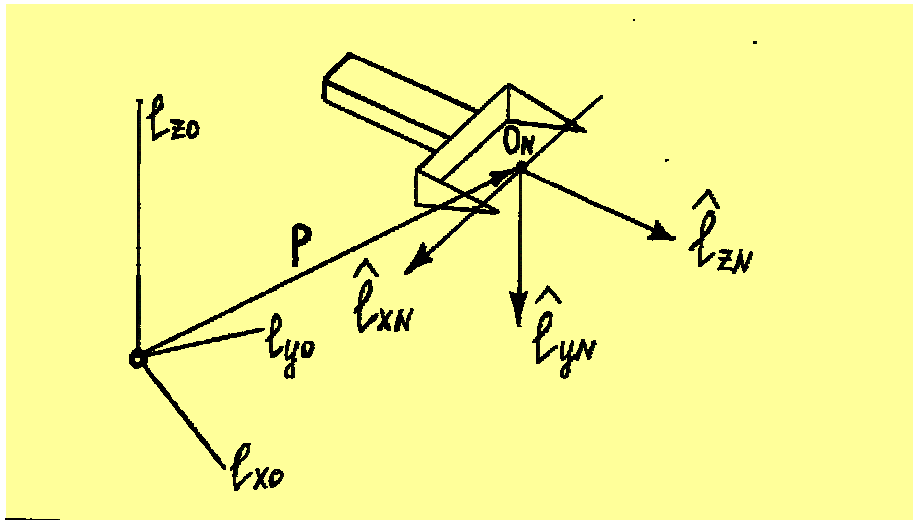


Рис. 1.7. Система координат манипулятора

Положение и ориентация схвата по отношению к абсолютной системе координат (0) определяется радиус-вектором \mathbf{P} и векторами $\hat{e}_{x_N}, \hat{e}_{y_N}, \hat{e}_{z_N}$, представленными в виде разложений по ортам системы (0).

Введем матрицу $T(q)$, определяющую правила пересчета координат из системы (\hat{N}) в систему (0).

$$T(q) = \begin{pmatrix} \hat{e}_{x_N} & \hat{e}_{y_N} & \hat{e}_{z_N} & \mathbf{P} \\ 0 & 0 & 0 & 1 \end{pmatrix}, \tag{1.30}$$

где в первой строке записаны вектор-столбцы размерности 3 в проекциях на оси системы (0).

Для того чтобы описать способ ее вычисления введем следующие обозначения.

Вектора, записываемые в виде разложения по системе координат (i) будем помечать верхним индексом (i), а по системе (\hat{i}) - индексом \hat{i} . Индекс (0) договоримся опускать. Вместо векторов размерности 3 часто будем использовать расширенные векторы, дописывая четвертую координату равную единице. В этих обозначениях пересчет координат из системы (\hat{N}) в абсолютную систему (0) в этих обозначениях будет выглядеть так

$$\begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} = T(q) \cdot \begin{pmatrix} x^{(\hat{N})} \\ y^{(\hat{N})} \\ z^{(\hat{N})} \\ 1 \end{pmatrix}, \tag{1.31}$$

где $q=(q_1, q_2, \dots, q_N)$.

Для описания пространственной структуры звеньев введем постоянные матрицы B_i , описывающие пересчет координат из системы (\hat{i}) в систему (i) на i-том звене

$$B_i = \begin{pmatrix} \hat{e}_{xi}^{(i)} & \hat{e}_{yi}^{(i)} & \hat{e}_{zi}^{(i)} & \hat{r}_i^{(i)} \\ 0 & 0 & 0 & 1 \end{pmatrix}, \tag{1.32}$$

где $\hat{r}_i^{(i)}$ - радиус – вектор в (i) – й системе координат, направленный в начало системы (\hat{i}). Матрицы B_i частично описывают геометрию звеньев. Относительное

смещение звеньев в i – й кинематической системе можно описать правилом пересчета при переходе от системы (i) к системе $(i - 1)$. Оно определяется матрицей

$$S_i(q_i) = \begin{pmatrix} e_{x_{i+1}}^{(i)} & e_{y_{i+1}}^{(i)} & e_{z_{i+1}}^{(i)} & r_{i+1}^{(i)} \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad (1.33)$$

зависящей от обобщенной координаты q_i . Матрицы $S_i(q_i)$ определяют способ соединения и относительного перемещения звеньев в кинематической паре.

Для любых пар поступательного типа

$$S_i(q_i) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & q_i \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad (1.34)$$

а для любых пар вращательного типа

$$S_i(q_i) = \begin{pmatrix} \cos q_i & -\sin q_i & 0 & 0 \\ \sin q_i & \cos q_i & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (1.35)$$

Матрица $T(q)$ перехода от системы координат (\hat{N}) к абсолютной системе координат (0) определяется следующим матричным произведением

$$T(q) = B_0 S_1(q_1) B_1 S_2(q_2) \dots S_N(q_N) B_N.$$

Элементы матрицы однозначно определяют положение и ориентацию в пространстве схвата манипулятора.

Приведем теперь математическую постановку задачи определения обобщенных координат манипулятора по заданному положению и ориентации схвата. Пусть требуется перевести схват в заданную точку $p^* = (x^*, y^*, z^*)^T$ и обеспечить нужную ориентацию системы координат (\hat{N}) , связанной со схватом. Эту ориентацию опишем ортогональной матрицей E^* с нормированными столбцами.

Составим из этих элементов матрицу T^* требуемого положения и ориентации схвата

$$T^* = \begin{pmatrix} E^* & P^* \\ 0 & 1 \end{pmatrix}, \quad (1.36)$$

где $0 = (0, 0, 0)$.

Введем функцию невязки

$$F(q) = \|T(q) - T^*\|_F, \quad (1.37)$$

в которой используется матричная норма Фробениуса, равная сумме квадратов элементов матрицы.

Поскольку существуют конструктивные ограничения, определяющие пределы изменения обобщенных координат манипулятора, то в задаче задается область изменения переменных q

$$D = \{q \in \mathbb{R}^N : a_i \leq q_i \leq b_i (i=1, \dots, N)\}. \quad (1.38)$$

Задача определения q^* , обеспечивающего наилучшее приближение в области D требуемого положения и ориентации схвата сводится к задаче оптимизации вида

$$F(q) \rightarrow \min, q \in D \in \mathbb{R}^N. \tag{1.39}$$

Сложный характер зависимости $F(q)$ определяет многоэкстремальность задачи. За счет возможной неоднозначности выбора q задача часто имеет несколько изолированных глобальных минимумов.

Более сложная постановка задачи возникает в том случае, когда имеются препятствия в зоне действия манипулятора. Пусть ограничения на пространственное положение манипулятора состоят в том, что центры всех его кинематических пар, находящихся в центрах систем (\hat{i}) , $(i=1, \dots, N)$ должны быть достаточно удалены от системы точек $V_j=(x_j, y_j, z_j)$, $(j=1, \dots, m)$, заданных в системе координат (0) .

Введем матрицы $T_i(q)$, описывающие положение и ориентацию в пространстве систем координат \hat{i} . Очевидно, что

$$T_i(q) = B_0 S_1(q_1) B_1 S_2(q_2) \dots S_i(q_i) B_i \tag{1.40}$$

Если обозначить через $P_i(q)$ радиус-вектор направленный из начала системы координат (0) в начало координат системы (\hat{i}) , являющейся центром i -той кинематической пары, то

$$\begin{pmatrix} P_i(q) \\ 1 \end{pmatrix} = T_i(q) \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}. \tag{1.41}$$

Введем функции расстояния $g_j(q)$ препятствия V_j от центров кинематических пар

$$g_j(q) = \max \left\{ \left\| \begin{pmatrix} V_j \\ 1 \end{pmatrix} - T_i(q) \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} \right\| : i = 1, \dots, N \right\}. \tag{1.42}$$

Дополнительные ограничения на обобщенные координаты состоят в том, что

$$q \in Q = \{q \in D : g_j(q) \geq \delta > 0 \ (j=1, \dots, m)\}. \tag{1.43}$$

В результате получим задачу с функциональными ограничениями

$$F(q) \rightarrow \min, q \in Q \in \mathbb{R}^N. \tag{1.44}$$

Выбор вспомогательных систем координат на звеньях и построение матриц $B_i, S_i(q_i)$ поясним на примере.

Пусть задан манипулятор, имеющий кинематическую схему, представленную на рис. 1.8.

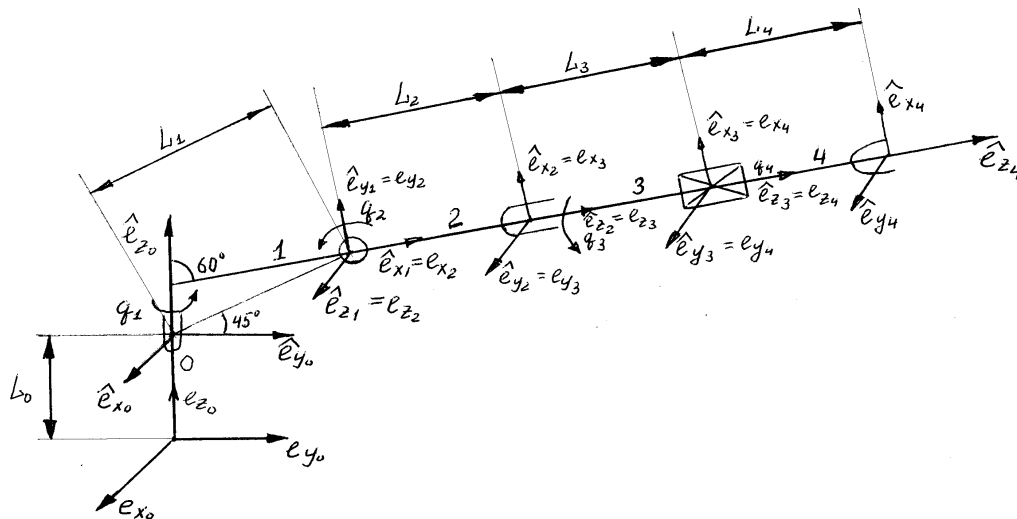


Рис. 1.8. Кинематическая схема манипулятора

Манипулятор показан в той конфигурации, которая принята за исходную. Она соответствует нулевым значениям обобщенных координат. Заметим, что в этой конфигурации положения в пространстве систем (\hat{i}) всегда совпадает с положением систем $(i+1)$.

Для этого манипулятора вид матриц, входящих в матричное описание $T(q)$ следующий

$$N=4, \quad q=(q_1, q_2, q_3, q_4)$$

$$B_0 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & L_0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad S_1(q_1) = \begin{pmatrix} \cos q_1 & -\sin q_1 & 0 & 0 \\ \sin q_1 & \cos q_1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (1.45)$$

$$B_1 = \begin{pmatrix} 0 & 0 & 1 & 0 \\ \frac{\sqrt{3}}{2} & -\frac{1}{2} & 0 & L_1 \frac{1}{2} \\ \frac{1}{2} & \frac{\sqrt{3}}{2} & 0 & L_1 \frac{1}{2} \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad S_2(q_2) = \begin{pmatrix} \cos q_2 & -\sin q_2 & 0 & 0 \\ \sin q_2 & \cos q_2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (1.46)$$

$$B_2 = \begin{pmatrix} 0 & 0 & 1 & L_2 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad S_3(q_3) = \begin{pmatrix} \cos q_3 & -\sin q_3 & 0 & 0 \\ \sin q_3 & \cos q_3 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (1.47)$$

$$B_3 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad S_4(q_4) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & q_4 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (1.48)$$

$$B_5 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & L_4 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (1.49)$$

Краткий обзор главы

В данной главе была изложена математическая постановка общей задачи оптимизации. Были введены основные понятия теории оптимизации, такие как вектор управляемых параметров, критерии эффективности, область поиска, область допустимых решений. Была введена классификация задач оптимизации, которая определяется

- количеством управляемых параметров (одномерная/ многомерная оптимизация);
- количеством критериев эффективности (однокритериальная/ многокритериальная оптимизация наличием функций ограничений (безусловная/ условная оптимизация) числом экстремумов у критериев эффективности (одноэкстремальная/ многоэкстремальная оптимизация);
- свойствами функций критериев эффективности и ограничений (линейность, непрерывность и т.д.);
- свойствами управляемых параметров (дискретность, непрерывность);
- свойствами области допустимых решений (односвязность, многосвязность).

Были рассмотрены два примера приложения теории оптимизации к решению практических задач.



Контрольные вопросы и упражнения

1. Дать определения основных понятий теории оптимизации (вектор варьируемых параметров, вектор критериев эффективности, вектор ограничений, область поиска экстремума, область допустимых решений).
2. Дать классификацию оптимизационных задач.
3. Привести общая схема численных методов многоэкстремальной оптимизации.
4. Сформулировать принципы создания учебно-исследовательских систем принятия решений.
5. Привести пример задачи оптимизации и указать ее классификацию.

Глава 2. Математические основы конструирования алгоритмов. Характеристические алгоритмы глобального поиска

2.1. Постановка задачи

Пусть $\varphi(x)$ - действительная функция, определенная в области Q N -мерного евклидова пространства R^N и принимающая во всех точках области конечные значения. Обозначим через

$$\inf_{x \in Q} \varphi(x) \equiv \inf\{\varphi(x) : x \in Q\} \quad (2.1)$$

точную нижнюю грань функции $\varphi(x)$ на множестве Q . Если существует точка x^* такая, что

$$\varphi(x^*) = \inf\{\varphi(x) : x \in Q\}, \quad (2.2)$$

то говорят, что функция $\varphi(x)$ достигает своей точной нижней грани на множестве Q , а точка называется *точкой глобального минимума* $\varphi(x)$ на множестве Q . Величина $\varphi(x^*)$ называется *наименьшим или глобально-минимальным значением* $\varphi(x)$ на множестве Q и обозначается

$$\min_{x \in Q} \varphi(x) \equiv \min\{\varphi(x) : x \in Q\} \quad (2.3)$$

Множество всех точек $x^* \in Q$, удовлетворяющих (2.2), будем обозначать через

$$Q^* = \mathop{\text{Arg min}}_{x \in Q} \varphi(x) \equiv \mathop{\text{Arg min}}\{\varphi(x) : x \in Q\} \quad (2.4)$$

Точка $\hat{x} \in Q$ называется *точкой локального минимума* функции на множестве Q , если существует такое число $\varepsilon > 0$, что для всех $x \in Q$ таких, что $\|x - \hat{x}\| < \varepsilon$, выполняется $\varphi(\hat{x}) \leq \varphi(x)$.

Определение 2.1. Задачей оптимизации будем называть задачу следующего вида:

найти заданные экстремальные характеристики функции $\varphi(x)$ на множестве Q .

Синонимически данную задачу также часто называют *задачей математического программирования*.

В зависимости от искомым экстремальных характеристик возможны различные постановки задачи оптимизации.

Постановка А. В качестве экстремальной характеристики рассматривается величина

$$\varphi^* = \inf\{\varphi(x) : x \in Q\} \quad (2.5)$$

Постановка В. Определить нижнюю грань из (2.5) и, если множество точек глобального минимума из (2.4) не пусто, найти хотя бы одну точку $x^* \in Q^*$.

Постановка С. Найти нижнюю грань φ^* из (2.5) и определить все точки глобального минимума (либо убедиться, что множество Q^* пусто).

Постановка D. Найти все точки и значения локальных минимумов.

Постановки А-D задачи математического программирования являются наиболее распространенными, хотя возможны и другие варианты, например, указанные в [1]. Искомые экстремальные характеристики, определяемые постановкой $V \in \{A, B, C, D\}$, назовем V -решением задачи математического программирования.

Символически общую задачу математического программирования будем записывать в виде

$$\varphi(x) \rightarrow \inf, x \in Q \quad (2.6)$$

и называть также *задачей минимизации* функции $\varphi(x)$ на множестве Q .

Так как точная верхняя грань функции $\varphi(x)$ на множестве Q

$$\sup\{\varphi(x) : x \in Q\} = -\inf\{-\varphi(x) : x \in Q\}, \quad (2.7)$$

то задача определения экстремальных характеристик, связанных с наибольшим значением функции $\varphi(x)$ (задача максимизации), сводится к задаче минимизации функции $-\varphi(x)$. Поэтому везде далее задача математического программирования будет рассматриваться в форме (2.6) и иногда называться просто задачей оптимизации. Функцию $\varphi(x)$ из (2.6) будем называть *целевой, минимизируемой* или *оптимизируемой функцией*, множество Q - *допустимой областью*, а элементы множества Q - *допустимыми точками*.

Задачу (2.6), для которой заведомо известно, что множество точек глобального минимума Q^* не пусто (достаточные условия непустоты Q^* даются, например, теоремой Вейерштрасса [2]), будем записывать как

$$\varphi(x) \rightarrow \min, x \in Q \quad (2.8)$$

Заметим, что часто задача математического программирования формулируется именно в таком виде.

2.2. Численные методы оптимизации

Сформулировав задачу минимизации, мы теперь должны дать ответ на основной вопрос: каким образом ее решать?

Классический подход математического анализа предлагает следующую процедуру аналитического решения задачи. Пусть $\varphi(x)$ - кусочно-гладкая на отрезке $[a,b]$ функция. Тогда минимум $\varphi(x)$ на $[a,b]$ может достигаться лишь в тех точках, где $\varphi'(x) = 0$, 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 $\varphi'(x) = 0$, 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 $\varphi(x)$ задается алгоритмически, т.е. в виде некоей расчетной схемы, когда по заданному аргументу x рассчитывается значение $\varphi(x)$. 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 (2.8) 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 (2.6) [3].

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 $\varphi(x)$) 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) [4]

$$s = \langle \{G_k\}, \{E_k\}, \{H_k\} \rangle \quad (2.9)$$

в котором

- $\{G_k\}$ - совокупность правил выбора точек испытаний, $k = 1, 2, \dots$;
- $\{E_k\}$ - совокупность правил построения приближенного решения (оценки экстремума), $k = 1, 2, \dots$;
- $\{H_k\}$ - совокупность правил остановки вычислительного процесса, $k = 1, 2, \dots$

Порядок проведения испытаний, или *вычислительная схема* алгоритма состоит в следующем.

1. Выбирается точка первого испытания

$$x^1 = G_1(\Phi) \in Q \quad (2.10)$$

2. Пусть выбрана точка k -го испытания $x^k \in Q$ ($k \geq 1$). Производится вычисление значения функции $z^k = \varphi(x^k)$. После этого имеем следующую поисковую (апостериорную) информацию о функции φ :

$$\omega_k = \{(x^1, z^1), (x^2, z^2), \dots, (x^k, z^k)\} \quad (2.11)$$

Полученная информация позволяет сузить класс, которому принадлежит функция $\varphi(x)$ до множества

$$\Phi(\omega_k) = \{\psi \in \Phi : \psi(x^i) = z^i, 1 \leq i \leq k\} \quad (2.12)$$

3. Определяется текущая оценка экстремума (приближенное решение)

$$e^k = E_k(\Phi, \omega_k) \quad (2.13)$$

4. Вычисляется точка очередного испытания

$$x^{k+1} = G_{k+1}(\Phi, \omega_k) \quad (2.14)$$

5. Определяется величина

$$h^k = H_k(\Phi, \omega_k) \in \{0, 1\}, \quad (2.15)$$

принимаящая 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 (2.8) 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

$$G_1(\Phi) = a, G_{k+1}(\Phi, \omega_k) = a + k \frac{b-a}{n}, k \geq 1 \quad (2.16)$$

$$H_k(\Phi, \omega_k) = \begin{cases} 0, & k = n+1 \\ 1, & k < n+1 \end{cases} \quad (2.17)$$

$$e^k = \varphi_k^*, \quad (2.18)$$

где

$$\varphi_k^* = \min_{1 \leq i \leq k} \varphi(x^i), \quad (2.19)$$

либо

$$e^k = (\varphi_k^*, x_k^*), \quad (2.20)$$

где

$$x_k^* = \arg \min_{1 \leq i \leq k} \varphi(x^i) \quad (2.21)$$

Контрольные вопросы и упражнения:

1. Предположим, что класс Φ представляет собой класс линейных функций и пусть проведено одно испытание в точке, обеспечившее результат z^1 . Что из себя будет представлять класс $\Phi(\omega_1)$?

2. Если затем проведено второе испытание $x^2 > x^1$ с результатом z^2 , как выглядит класс $\Phi(\omega_2)$?

3. Пусть для той же задачи после третьего испытания в точке $x^1 < x^3 < x^2$ получено значение z^3 такое, что $z^3 < z^1$, $z^3 < z^2$, какие выводы можно сделать?

Итак, решая задачу минимизации функции $\varphi(x)$, метод поиска порождает (сопоставляет функции) последовательность $x^1, x^2, \dots, x^k, \dots$ координат испытаний, или просто последовательность испытаний (x^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. $z^i = \varphi(x^i)$). При

этом свойства метода определяются свойствами последовательностей $\{x^k\}$ 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 $\{x^k\}$ 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 $\{x^k\}$. 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 (2.9) $H_k = 1$ для любого k , т.е. условие остановки отсутствует) и требуя, чтобы эта последовательность сходилась к точному решению задачи. Поскольку в постановках В-D искомое решение может содержать *несколько* точек минимума, сходимость метода будем понимать в смысле следующего определения.

Определение 2.2. Последовательность испытаний $\{x^k\}$ сходится к решению задачи (2.6), определенному соответствующей постановкой, если:

- 1) она содержит подпоследовательность $\{\bar{x}^k\}$, для которой

$$\lim_{k \rightarrow \infty} \varphi(\bar{x}^k) = \varphi^* ;$$
- 2) в случае, когда решение включает одну или несколько точек минимума, для каждой такой точки существует сходящаяся к ней подпоследовательность последовательности $\{x^k\}$.

Последовательность испытаний, сходящаяся к точному решению постановки $V \in \{A, B, C, D\}$, будем называть *минимизирующей* последовательностью для постановки V , либо V -минимизирующей последовательностью. Термин "минимизирующая последовательность" введен в [2] и соответствует понятию A -минимизирующей последовательности.

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 функций, т.е. априорно известно, что минимизируемая функция $\varphi(x)$ непрерывна в области Q . Предположим, что вычислены значения функции $\varphi(x)$ в конечном числе точек x^1, x^2, \dots, x^k . Что после этого можно сказать о координате глобального минимума? Каковы бы ни были точки x^1, x^2, \dots, x^k и значения z^1, z^2, \dots, z^k , для любой точки $x^* \in Q$ ($x^* \notin \{x^1, \dots, x^k\}$) всегда можно построить непрерывную функцию, проходящую через точки $(x^i, z^i), 1 \leq i \leq k$, т.е. принадлежащую классу $\Phi(\omega_k)$ из (2.12),

которая имеет глобальный минимум в точке x^* с любым наперед заданным значением $\varphi^* < \min_{1 \leq i \leq k} z^i$.

$$1 \leq i \leq k$$

Например, в качестве такой функции можно взять интерполяционный полином k -й степени, проходящий через точки $(x^i, z^i), 1 \leq i \leq k$, и точку (x^*, φ^*) .

Все сказанное означает, что по результатам конечного числа испытаний никаких выводов о расположении координаты глобального минимума сделать нельзя. Точно так же о величине φ^* глобального минимума можно лишь сказать, что

$$\varphi^* \leq \varphi_k^*, \quad (2.22)$$

где φ_k^* из (2.19), однако оценить величину

$$\varepsilon_k = |\varphi^* - \varphi_k^*|, \quad (2.23)$$

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 $\varphi(x)$.

Для примера рассмотрим класс *строго унимодальных* на отрезке $[a, b]$ функций, т.е. функций, для каждой из которых существует точка $x^* \in [a, b]$ такая, что на отрезке $[a, x^*]$ функция строго убывает, а на отрезке $[x^*, b]$ - строго возрастает. Пусть проведены испытания в точках x^1 и x^2 интервала (a, b) и получены значения целевой функции z^1 и z^2 . Предположим, что $x^1 < x^2$ и $z^1 < z^2$. Тогда в силу унимодальности очевидно, что на отрезке $[x^2, b]$ точка минимума x^* находится не может, и в качестве области локализации координаты минимума можно рассмотреть полуинтервал $[a, x^2)$, называемый *интервалом неопределенности*.

Пусть теперь в общем случае проведено k испытаний в точках $x^1, x^2, \dots, x^k \in (a, b)$ 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.

$$a = x_0 < x_1 < x_2 < \dots < x_k < x_{k+1} = b \quad (2.24)$$

Тогда интервалом неопределенности будет интервал (x_{i-1}, x_{i+1}) , где номер i определяется из условия $x_i = x_k^*$, где x_k^* из (2.21) (в случаях $i=1$ и $i=k$ интервалами неопределенности будут полуинтервалы $[a, x_2)$ и $(x_{k-1}, b]$ 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) длиной этого интервала, ибо

$$|x_k^* - x^*| < \varepsilon = x_{i+1} - x_{i-1} \quad (2.25)$$

Что касается величины глобального минимума, то строгой унимодальности для получения оценки (2.23) 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. Детальные схемы такого оценивания могут быть найдены в работах [3-6].

2.3. Оптимальность 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 [4] затраты на поиск с некоторой мерой близости оценки e^k из (2.13) к точному решению задачи.

Приведем формальную постановку задачи определения эффективности алгоритма оптимизации [3, 4, 7]. Согласно этой постановке рассматривается некоторый класс S алгоритмов $s \in S$, предназначенных для решения задач минимизации (2.8) функций φ из класса Φ . Вводится вещественная функция $L(\varphi, s)$, называемая критерием эффективности, которая количественно характеризует эффективность решения задачи минимизации функции $\varphi \in \Phi$ с помощью метода $s \in S$. Для определенности будем считать, что чем меньше величина $L(\varphi, s)$, тем выше эффективность алгоритма.

Следуя общей схеме теории исследования операций [8], в качестве эффективности метода s на классе Φ примем либо гарантированный результат

$$W(s) = \sup_{\varphi \in \Phi} L(\varphi, s) \quad (2.26)$$

(наихудшее возможное "достижение" алгоритма на классе Φ), либо среднюю по классу Φ эффективность

$$W(s) = \int_{\Phi} L(\varphi, s) dP(\varphi), \quad (2.27)$$

где $P(\varphi)$ - 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

$$W(s^*) = \inf_{s \in S} W(s) \quad (2.28)$$

Если точная нижняя грань в (2.28) не реализуема (s^* не существует), можно рассмотреть ε -оптимальный алгоритм s_ε^* , для которого (при фиксированном $\varepsilon > 0$) имеет место

$$W(s_\varepsilon^*) \leq \inf_{s \in S} W(s) + \varepsilon \quad (2.29)$$

Если следовать терминологии исследования операций, то методы оптимизации в ситуации гарантированного результата (2.26) можно назвать также стратегиями оптимизации, а оптимальные (ε -оптимальные) алгоритмы – оптимальными (ε -оптимальными) минимаксными стратегиями. При рассмотрении критерия (2.27), который в теории статистических решений называют функцией риска, стратегия, минимизирующая риск, называется байесовской, что приводит к использованию термина "*байесовские методы оптимизации*".

Формулировка критериев оптимальности в виде (2.26), (2.27) 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 информации о функции. Это приводит для критерия (2.26) к ориентации на худший, а для критерия (2.27) – на "типичный" (наиболее вероятный) случай, хотя оптимизируемая функция может существенно отличаться от худшей, либо типичной. В связи с этим А.Г.Сухарев ввел понятие *последовательно-оптимального* (в терминологии [9] – *наилучшего*) 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 [3]). Понятие последовательной оптимальности можно применить и к байесовским методам поиска [4].

Проблема отыскания оптимального алгоритма решения экстремальной задачи (2.6) в свою очередь сводится к решению экстремальной задачи (2.28), причем, как правило, существенно более сложной, поскольку областями поиска являются нечисловые множества. Как указано в [4], решение этой задачи возможно лишь при наличии соответствующего математического аппарата исследования функции $L(\varphi, s)$, которая, будучи мерой эффективности решения задачи (2.6), 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 унимодальных функций одной переменной. Именно для этого класса в работе [10] 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 [2, 3, 11, 12].

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 [3, 13-15]. Оптимальные минимаксные алгоритмы построены в работах [3, 8, 14, 16-20]. При этом оказалось, что вследствие ориентации на худший случай, которым является функция-константа, многие оптимальные алгоритмы представляют собой метод перебора по равномерной сетке. Более экономные последовательно-оптимальные алгоритмы удалось построить А.Г.Сухареву [9] для некоторых подклассов липшицевых функций.

Что касается байесовых методов поиска экстремума, то Й.Б.Моцкусом [21] показано, что задача их построения может быть сведена к решению некоторой системы

рекуррентных уравнений. Однако и в этом случае задача остается настолько сложной, что ни одного последовательно оптимального метода получить не удалось, а оптимальные байесовы методы, построенные в работах [22, 23], точно не реализуемы.

Трудности конструирования оптимальных в смысле (2.28) алгоритмов поиска привели к использованию более простых понятий оптимальности. Одним из таких понятий является так называемая *одношаговая оптимальность* [3, 4, 5, 7], когда алгоритм размещает очередное испытание наилучшим образом, предполагая, что оно является последним.

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 $\varphi \in \Phi_k$, меняющееся от шага к шагу. В этом случае использование принципа одношаговой оптимальности совершенно естественно.

Приведем формальное описание *одношагово-оптимального метода поиска экстремума* [4, 5]. С этой целью вернемся к формальной схеме метода оптимизации (2.9) и для любого возможного ω_k из (2.11) определим класс $\Phi(\omega_k)$ из (2.12), 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 $\varphi \in \Phi(\omega_k)$ после первых k испытаний реализует данное ω_k . Введем последовательность функций $L_k(\varphi, \omega_k)$, определяющих эффективность текущей оценки экстремума e^k для функции $\varphi \in \Phi(\omega_k)$, которую за k испытаний обеспечил алгоритм $s \in S(\omega_k)$. Тогда алгоритм s может быть описан набором

$$s = \langle \{G_k\}, \{E_k\}, \{H_k\}, \{L_k\} \rangle \quad (2.30)$$

Множество таких алгоритмов обозначим через S_0 . Введем величину

$$W_{k+1}(x^{k+1}) = \sup_{\varphi \in \Phi(\omega_k)} L_{k+1}(\varphi, \omega_{k+1}) \quad (2.31)$$

Тогда метод $s \in S_0$ называется (минимаксным) одношагово-оптимальным, если точки испытаний x^k , $k = 1, 2, \dots$, порождаемые методом s , удовлетворяют условиям

$$W_k(x^k) = \min_{\tilde{x}^k \in Q} W_k(\tilde{x}^k) \quad (2.32)$$

Если вместо гарантированного результата (2.31) рассмотреть математическое ожидание

$$W_{k+1}(x^{k+1}) = \int_{\Phi(\omega_k)} L_{k+1}(\varphi, \omega_{k+1}) dP(\omega_k), \quad (2.33)$$

где $P(\omega_k)$ - условное по отношению к результатам испытаний распределение

вероятностей, то метод, определяемый соотношением (2.32), будет называться байесовским одношагово-оптимальным алгоритмом.

Минимаксные одношагово-оптимальные алгоритмы предложены в работах [6, 24, 25], а при байесовском подходе одношагово-оптимальные методы построены в [4,5, 26-30].

Дальнейшим упрощением принципа оптимальности является рассмотрение *асимптотически оптимальных* алгоритмов. Для определения этого понятия обозначим через S_N множество методов (2.9), в которых остановка осуществляется ровно через N шагов поиска. Введем величину

$$W(N) = \inf_{s \in S_N} W(s), \quad (2.34)$$

где $W(s)$ из (2.26) или (2.27). Для оптимального метода $s_N^* \in S_N$ очевидно выполняется $W(s_N^*) = W(N)$.

Алгоритм \hat{s}_N называется *асимптотически оптимальным*, если $W(\hat{s}_N) / W(N) \rightarrow 1$ при $N \rightarrow \infty$. Обсуждение различных вопросов, связанных с указанным понятием оптимальности, можно найти в [11].

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 можно поставить оптимальность согласно (2.28) (или ε -оптимальность (2.29)), ориентированную только на априорное знание. Далее следует принцип оптимальности на один шаг вперед и, наконец, асимптотическая оптимальность.

2.4. 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 (2.6). 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, называемого классом характеристических алгоритмов и включающего многие известные алгоритмы, созданные в рамках различных подходов к конструированию методов оптимизации.

Рассмотрим одномерную задачу (2.6) для области поиска $Q = [a, b]$ и класс методов оптимизации (2.9), 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 $\{x^k\} = x^1, x^2, \dots, x^k, \dots$, изучение свойств которой и будет предметом нашего интереса.

Определение 2.3. Алгоритм решения задачи (2.6) называется *характеристическим*, если, начиная с некоторого шага поиска $k_0 \geq 1$, выбор

координаты x^{k+1} очередного испытания ($k \geq k_0$) заключается в выполнении следующих действий.

- 1) Задать набор

$$\Lambda_k = \{x_0, x_1, \dots, x_\tau\} \quad (2.35)$$

конечного числа $\tau+1 = \tau(k)+1$ точек области $Q = [a, b]$, полагая, что $a \in \Lambda_k, b \in \Lambda_k$, все координаты предшествующих испытаний $x^i \in \Lambda_k, 1 \leq i \leq k$, и множество Λ_k упорядочено (нижним индексом) по возрастанию координаты, т.е.

$$a = x_0 < x_1 < \dots < x_{\tau-1} < x_\tau = b. \quad (2.36)$$

- 2) Каждому интервалу (x_{i-1}, x_i) , $1 \leq i \leq \tau$, поставить в соответствие число $R(i)$, называемое характеристикой этого интервала.
- 3) Определить интервал (x_{t-1}, x_t) , которому соответствует максимальная характеристика $R(t)$, т.е.

$$R(t) = \max \{R(i) : 1 \leq i \leq \tau\} \quad (2.37)$$

- 4) Провести очередное испытание в точке

$$x^{k+1} = D(t) \in (x_{t-1}, x_t). \quad (2.38)$$

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 (x_{i-1}, x_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 (2.35) наряду с координатами испытаний может содержать точки, в которых испытания не проводились (например, в ряде информационно-статистических алгоритмов [5] 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 (2.36). Так, координата i -го испытания x^i в множестве Λ_k получит нижний индекс j , т.е. $x^i = x_j$, причем от шага к шагу номер j может меняться ($j = j(k)$).

Понятие характеристичности метода оптимизации впервые было введено В.А.Гришагиным [31] и позднее обобщено и распространено на другие классы задач и типы алгоритмов [32-34, 39].

В качестве иллюстрации приведем примеры известных алгоритмов глобальной оптимизации. В этих алгоритмах два первых испытания проводятся в точках $x^1 = a$ и $x^2 = b$, характеристическое правило вступает в действие, начиная с $k = 2$, при этом множество Λ_k ($k \geq 2$) состоит только из точек испытаний, т.е. $\Lambda_k = \{x^1, x^2, \dots, x^k\}$ и, следовательно, $\tau = k - 1$. Будем также использовать обозначение $z_j = \varphi(x_j)$ для значений целевой функции в точках $x_j \in \Lambda_k$.

Метод последовательного сканирования (перебор).

Для этого метода характеристикой интервала является его длина, т.е.

$$R(i) = x_i - x_{i-1}, \quad (2.39)$$

а точка очередного испытания выбирается в середине самого длинного интервала:

$$x^{k+1} = 0.5(x_{t-1} + x_t). \quad (2.40)$$

Метод ломаных.

В данном методе, который был построен С.А.Пиявским [6] для оптимизации липшицевых функций, характеристика

$$R(i) = 0.5m(x_i - x_{i-1}) - (z_i + z_{i-1})/2, \quad (2.41)$$

а точка очередного испытания выбирается согласно выражению

$$x^{k+1} = 0.5(x_t + x_{t-1}) - (z_t - z_{t-1})/(2m), \quad (2.42)$$

где $m > 0$ - параметр метода.

Информационно-статистический алгоритм глобального поиска (АГП).

Обсуждаемый метод предложен Р.Г.Стронгиным [4] как байесовский одношагово-оптимальный алгоритм и использует характеристику

$$R(i) = m(x_i - x_{i-1}) + \frac{(z_i - z_{i-1})^2}{m(x_i - x_{i-1})} - 2(z_i + z_{i-1}), \quad (2.43)$$

а точку нового испытания формирует согласно (2.42). Величина $m > 0$ вычисляется в соответствии с выражением

$$m = \begin{cases} rM, & M > 0 \\ 1, & M = 0 \end{cases}, \quad (2.44)$$

где

$$M = \max_{1 \leq i \leq \tau} \frac{|z_i - z_{i-1}|}{x_i - x_{i-1}}, \quad (2.45)$$

а $r > 1$ - параметр метода.

Контрольные вопросы и упражнения:

1. Выполните несколько первых итераций метода сканирования при минимизации функции $\varphi(x)$ 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 $\varphi(x) = x$ 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 (2.44) при $r = 2$. Попытайтесь 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 2.1. Одношагово-оптимальный (минимаксный или байесовский) алгоритм является характеристическим.

Доказательство. Перепишем условие одношаговой оптимальности (2.32) в виде

$$W_k(x^k) = \min_{\tilde{x}^k \in [a, b]} W_k(\tilde{x}^k) = \min_{1 \leq i \leq \tau} \min_{\tilde{x}^k \in [x_{i-1}, x_i]} W_k(\tilde{x}^k) \quad (2.46)$$

Отсюда следует, что в качестве характеристики интервала (x_{i-1}, x_i) можно взять величину

$$R(i) = - \min_{\tilde{x}^k \in [x_{i-1}, x_i]} W_k(\tilde{x}^k), \quad (2.47)$$

а

$$D(t) = \arg \min_{\tilde{x}^k \in [x_{t-1}, x_t]} W_k(\tilde{x}^k)$$

Теорема доказана.

Теорема 2.2. Пусть точка x^* является предельной точкой (точкой накопления) последовательности поисковых испытаний $\{x^k\}$, порождаемой характеристическим алгоритмом при решении задачи (2.6) на отрезке $[a, b]$, причем $x^* \neq a$ и $x^* \neq b$.

Предположим, что характеристики $R(i)$ и правила выбора точки очередного испытания $D(t)$ удовлетворяют следующим требованиям:

а) если при $k \rightarrow \infty$ точка $\bar{x} \in [x_{i(k)-1}, x_{i(k)}]$ и $x_{i(k)-1} \rightarrow \bar{x}$, $x_{i(k)} \rightarrow \bar{x}$, тогда

$$R(i(k)) \rightarrow -\mu\varphi(\bar{x}) + c \quad (2.48)$$

б) в случае, когда, начиная с некоторого шага поиска, интервал (x_{i-1}, x_i) , $i = i(k)$ не содержит точек поисковых испытаний, т.е. существует $\tilde{k} \geq 1$ такой, что для всех $k \geq \tilde{k}$

$$(x_{i-1}, x_i) \cap \{x^k\} = \emptyset, \quad (2.49)$$

для характеристики интервала справедливо

$$\lim_{k \rightarrow \infty} R(i) > -\mu \min\{\varphi(x_{i-1}), \varphi(x_i)\} + c; \quad (2.50)$$

$$c) \max\{x^{k+1} - x_{i-1}, x_i - x^{k+1}\} \leq \nu(x_i - x_{i-1}), \quad (2.51)$$

где μ, c, ν - некоторые константы, причем $\mu \geq 0$, $0 < \nu < 1$.

Тогда последовательность $\{x^k\}$ содержит две подпоследовательности, одна из которых сходится к x^* слева, а другая справа.

Доказательство. Рассмотрим вначале случай, когда $x^* \notin \{x^k\}$. Обозначим через $p = p(k)$, $k \geq 1$, номер интервала (x_{p-1}, x_p) , содержащего на k -м шаге поиска предельную точку x^* . Очевидно, что для $k=1$ $[x_{p-1}, x_p] = [a, b]$. После попадания в интервал (x_{p-1}, x_p) точки очередного испытания x^{k+1} (в этом случае $p=t$) для нового интервала $(x_{p(k+1)-1}, x_{p(k+1)})$, содержащего x^* , согласно (2.51) справедлива оценка

$$x_{p(k+1)-1} - x_{p(k+1)} \leq \nu(x_{p(k)-1} - x_{p(k)}).$$

Но тогда после попадания с начала поиска s испытаний в интервал с точкой x^* его длина будет удовлетворять неравенству

$$x_{p-1} - x_p \leq \nu^s (b - a). \quad (2.52)$$

Поскольку точка x^* предельная, после образования на некотором шаге интервала (x_{p-1}, x_p) в него попадет бесконечное число испытаний, поэтому из (2.52) следует, что

$$\lim_{k \rightarrow \infty} (x_{p(k)-1} - x_{p(k)}) = 0, \quad (2.53)$$

На основании (2.53) в качестве искомым подпоследовательностей мы можем взять последовательность $\{x_{p(k)-1}\}$ левых и последовательность $\{x_{p(k)}\}$ правых концов интервалов, содержащих x^* .

Пусть теперь найдется номер $q \geq 1$ такой, что $x^q = x^*$. Тогда при любом $k \geq q$ существует номер $j = j(k)$, $0 \leq j \leq \tau$, для которого $x_j = x^*$. Допустим, что имеет место односторонняя сходимости к x^* , например, слева. Тогда найдется номер $\tilde{k} \geq q$ такой, что при $k \geq \tilde{k}$ испытания в интервал (x_j, x_{j+1}) попадать не будут.

Из (2.49), (2.50) для интервала (x_j, x_{j+1}) и как следствие соотношения (2.48) для интервала (x_{j-1}, x_j) мы получаем, что

$$\lim_{k \rightarrow \infty} R(j+1) > -\mu \min\{\varphi(x_j), \varphi(x_{j+1})\} + c \geq -\mu\varphi(x^*) + c$$

$$\lim_{k \rightarrow \infty} R(j) = -\mu\varphi(x^*) + c$$

откуда, начиная с некоторого шага поиска, будет следовать выполнимость неравенства

$$R(j+1) > R(j) \tag{2.54}$$

Однако вследствие решающего правила (2.35)-(2.38) соотношение (2.54) противоречит невозможности проведения испытаний в интервале (x_j, x_{j+1}) . Данное противоречие завершает доказательство.

Следствие. В вычислительную схему характеристического алгоритма, удовлетворяющую условиям Теоремы 2.2, можно ввести условие остановки вида

$$x_t - x_{t-1} \leq \varepsilon, \tag{2.55}$$

где t из (2.37), 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 x^* . Обозначим через $p = p(k)$, $k \geq 1$, номер интервала, содержащего точку x^* на k -м шаге поиска. Т.к. данная точка предельная, то в интервал (x_{p-1}, x_p) попадет бесконечное число испытаний и для него будет иметь место соотношение (2.53), из которого следует, что условие (2.55) неизбежно выполнится на некотором шаге поиска. Заметим, что если точка x^* не является внутренней, для справедливости (2.53) достаточно односторонней сходимости.

Проверим выполнимость условий Теоремы 2.2 для рассмотренных примеров характеристических алгоритмов.

Метод последовательного сканирования.

Если интервал (x_{i-1}, x_i) стягивать в точку, характеристика метода (2.39) стремится к нулю. Поэтому в (2.48) в качестве констант μ и c можно взять $\mu=c=0$. Если же в интервал (x_{i-1}, x_i) , начиная с некоторого шага поиска, испытания попадать не будут, то его длина (совпадающая с характеристикой), будет оставаться положительной, что обеспечит выполнимость условия (2.50) для выбранных μ и c . Что касается неравенства (2.51), то оно очевидно выполняется при $\nu = 0.5$

Метод ломаных.

Для рассмотрения условий теоремы сделаем предположение, что минимизируемая функция $\varphi(x)$ удовлетворяет условию Липшица с константой $L > 0$, т.е.

$$|\varphi(x') - \varphi(x'')| \leq L|x' - x''|, x', x'' \in [a, b], \quad (2.56)$$

и, кроме того, параметр метода $m > L$.

В силу липшицевости функция $\varphi(x)$ непрерывна, и, следовательно, при стягивании интервала (x_{i-1}, x_i) к точке \bar{x} характеристика интервала будет стремиться к величине $-\varphi(\bar{x})$, т.е. для выполнимости (2.48) можно положить $\mu=1$ и $c=0$. Проверим теперь для данных μ и c справедливость неравенства (2.50) для интервала (x_{i-1}, x_i) , удовлетворяющего (2.50). Воспользуемся простым соотношением

$$\min\{z_{i-1}, z_i\} = \frac{1}{2}(z_{i-1} + z_i - |z_{i-1} - z_i|) \quad (2.57)$$

и оценим характеристику (2.41) метода, воспользовавшись условием Липшица и учитывая, что длина интервала остается, начиная с некоторого шага поиска, положительной и неизменной:

$$\begin{aligned} 0.5m(x_i - x_{i-1}) - (z_i + z_{i-1})/2 &> 0.5L(x_i - x_{i-1}) - (z_i + z_{i-1})/2 \geq 0.5(z_{i-1} + z_i - |z_{i-1} - z_i|) = \\ &= -\min\{z_{i-1}, z_i\} \end{aligned}$$

Полученные соотношения устанавливают справедливость (2.50).

Для определения величины ν в (2.51) оценим величину

$$\begin{aligned} x_i - x^{k+1} &= 0.5(x_i - x_{i-1}) + (z_i - z_{i-1})/(2m) \leq 0.5(x_i - x_{i-1}) + L(x_i - x_{i-1})/(2m) = \\ &= 0.5(1 + L/m)(x_i - x_{i-1}) \end{aligned}$$

Аналогичная оценка имеет место и для интервала (x_{i-1}, x^{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 (2.56) целевой функции, нетрудно показать, что данный метод также удовлетворяет условиям теоремы 2.2 с $\mu=4$, $c=0$ и $\nu=0.5(1+1/r)$.

Вследствие липшицевости при стягивании интервала (x_{i-1}, x_i) к точке \bar{x} его характеристика будет стремиться к величине $-4\varphi(\bar{x})$, т.е. для выполнимости (2.48) можно положить $\mu=4$ и $c=0$.

Для проверки (2.50) при условии (2.49) рассмотрим два случая. Пусть сначала для интервала (x_{i-1}, x_i) справедливо $z_{i-1} = z_i$. Тогда характеристика

$$R(i) = m(x_i - x_{i-1}) - 2(z_i + z_{i-1}) > -2(z_i + z_{i-1}) = -4 \min\{z_{i-1}, z_i\},$$

поскольку длина интервала (x_{i-1}, x_i) , начиная с некоторого шага поиска k_Δ , перестает изменяться, т.е. существует константа $\Delta > 0$ такая, что при $k > k_\Delta$ имеет место $x_i - x_{i-1} > \Delta > 0$.

Предположим теперь, что $z_{i-1} \neq z_i$. Представим характеристику (2.43) в виде

$$R(i) = |z_i - z_{i-1}| \left(\beta + \frac{1}{\beta} \right) - 2(z_i + z_{i-1}),$$

где вследствие (2.44), (2.45) $0 < \beta = \frac{m(x_i - x_{i-1})}{|z_i - z_{i-1}|} < 1$. В этом случае величина

$\beta + \frac{1}{\beta} > 2$, 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 (2.51), то в соответствии с (2.45) справедлива оценка

$$\begin{aligned} x_t - x^{k+1} &= 0.5(x_t - x_{t-1}) + (z_t - z_{t-1})/(2m) \leq 0.5(x_t - x_{t-1}) + M(x_t - x_{t-1})/(2m) = \\ &= 0.5(1+1/r)(x_t - x_{t-1}) \end{aligned}$$

Длину интервала (x_{t-1}, x^{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 2.3. Если в условиях теоремы 2.2 $\mu=0$, тогда любая точка области поиска является предельной точкой последовательности поисковых испытаний.

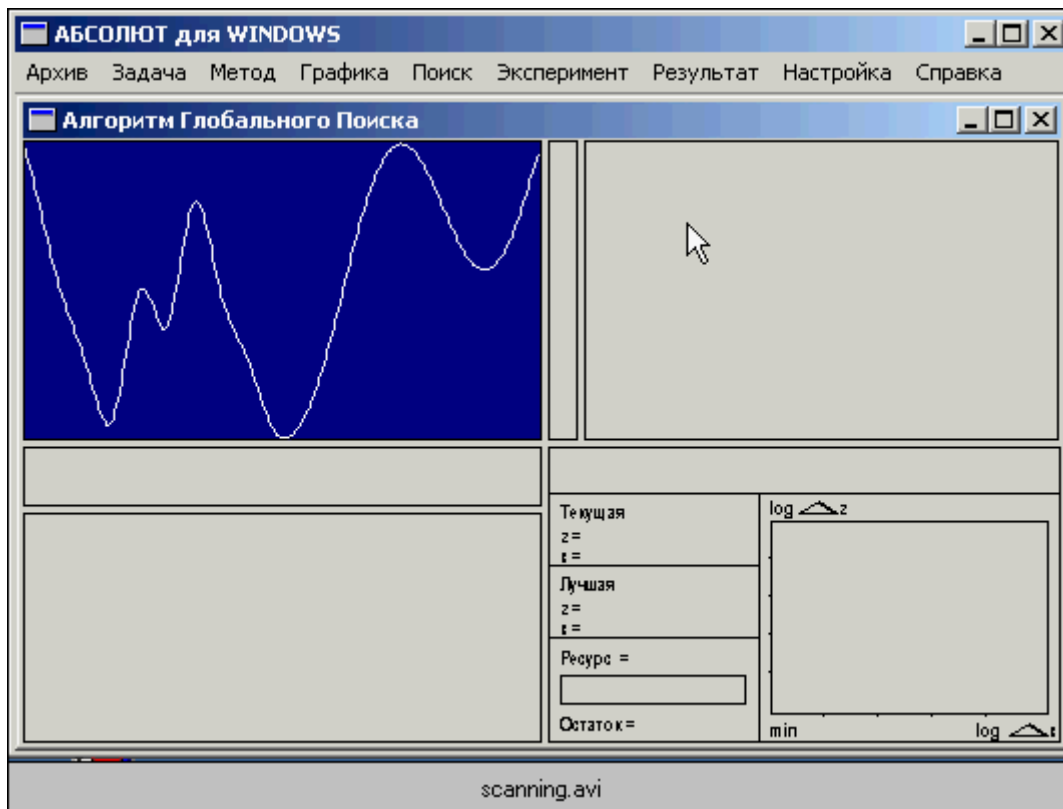
Доказательство. Предположим, что некоторая точка x' не является предельной для последовательности поисковых испытаний, порождаемой алгоритмом. Это означает, что начиная с некоторого шага поиска, испытания в интервал (x_{q-1}, x_q) , $q = q(k)$, попадать не будут, и тогда согласно (2.50) характеристика $R(q) > 0$. Но

последовательность испытаний, будучи ограниченной пределами отрезка $[a, b]$, содержит хотя бы одну сходящуюся подпоследовательность. Для интервала (x_{p-1}, x_p) , $p = p(k)$, содержащего предельную точку данной подпоследовательности, вследствие двусторонней сходимости справедливо соотношение (2.48), 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 (2.37) 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. Среди примеров, которые мы рассмотрели, подобной сходимостью обладает метод перебора, а среди других известных алгоритмов – методы [26, 28, 29]. 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.

Для иллюстрации поведения методов данного типа продемонстрируем динамику поиска метода перебора. Под графиком минимизируемой функции штрихами изображены координаты проведенных испытаний.



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 [39].

Другой тип поведения характеристических алгоритмов устанавливает

Теорема 2.4. Пусть в условиях теоремы 2.2 $\mu > 0$ и x^* - предельная точка последовательности поисковых испытаний. Тогда

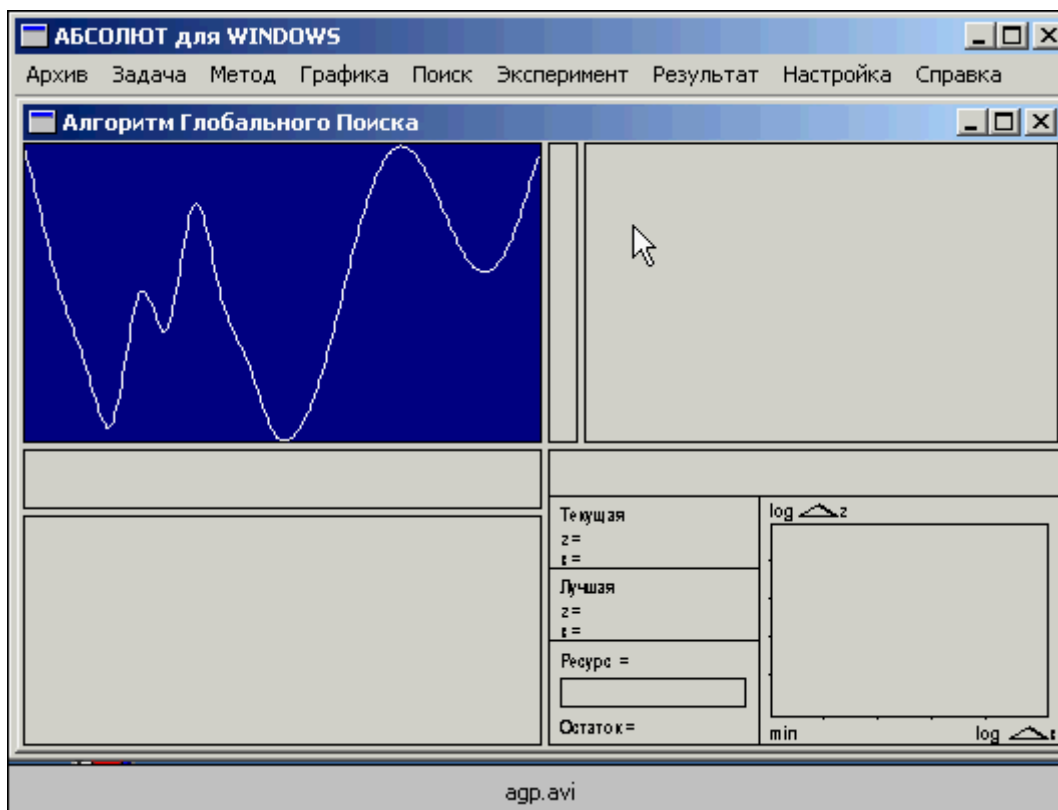
- 1) $\varphi(x^k) \geq \varphi(x^*)$, $k \geq 1$;
- 2) если существует еще одна предельная точка $x^{**} \neq x^*$, то $\varphi(x^{**}) = \varphi(x^*)$;
- 3) если в области поиска функция $\varphi(x)$ имеет конечное число локальных минимумов, то x^* является точкой локального минимума целевой функции в области Q .

Доказательство может быть найдено в [33].

Итак, если в условиях (2.48), (2.50) 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 [4,5, 35, 36], а также методы [32, 40].

Для примера проиллюстрируем работу АГП. Как и ранее, штрихами под графиком функции отмечаются координаты испытаний метода.



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 2.4 не гарантирует сходимость к глобальному минимуму исследуемой задачи оптимизации. Такие гарантии (достаточные условия сходимости) дает

Теорема 2.5. Пусть функция $\varphi(x)$ удовлетворяет на отрезке $[a, b]$ условию Липшица (2.56). Тогда точка x^* глобального минимума функции $\varphi(x)$ на данном отрезке является предельной точкой последовательности поисковых испытаний, порождаемой характеристическим алгоритмом при решении задачи (2.6), если выполняются условия (2.48) и (2.51) теоремы 2.2, а для интервала (x_{i-1}, x_i) со свойством (2.50) справедливо неравенство

$$\lim_{k \rightarrow \infty} R(i) > \frac{\mu}{2}(L(x_i - x_{i-1}) - \varphi(x_{i-1}) - \varphi(x_i)) + c; \quad (2.58)$$

где константы μ, c и ν удовлетворяют условиям теоремы 2.2.

Доказательство теоремы приведено в [33].

Следствие 1. При $\mu = 0$ (2.58) совпадает с (2.50) и теорема 2.5 становится идентичной теореме 2.3; при этом липшицевость целевой функции необязательна.

Следствие 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 2.4.

Заметим, что для метода ломаных (2.58) очевидно выполняется при $m > L$, а для АГП – при $m > 2L$.

2.5. Индексная схема учета ограничений

Рассмотрим одномерную задачу (2.8) для случая, когда область Q задается с помощью ограничений-неравенств, т.е.

$$Q = \{x \in [a, b], g_i(x) \leq 0, 1 \leq i \leq m\} \quad (2.59)$$

и перепишем данную задачу в форме

$$\min \{\varphi(x) : x \in [a, b], g_i(x) \leq 0, 1 \leq i \leq m\} \quad (2.60)$$

Задаче (2.60) можно сопоставить вспомогательную задачу, которая имеет решение даже в случае, когда допустимая область является пустой. Для этого построим классификацию точек $x \in [a, b]$ по числу $\nu = \nu(x)$ выполняющихся в них ограничений. При этом индекс $\nu = \nu(x)$ определяется условиями

$$x \in Q_\nu, x \notin Q_{\nu+1}, \tag{2.61}$$

где

$$Q_1 = [a, b], Q_{i+1} = \{x \in Q_i : g_i(x) \leq 0\}, 1 \leq i \leq m\}. \tag{2.62}$$

Обозначим через M максимальное значение индекса в области поиска и введем вспомогательную задачу

$$g_M^* = g_M(x^*) = \min\{g_M(x) : x \in Q_M\}, \tag{2.63}$$

Область Q_M не пуста и, следовательно, задача (2.63) всегда имеет решение.

Для решения задачи (2.60) предложен [5, 37, 38] следующий индексный алгоритм, который может быть реализован в рамках характеристической вычислительной схемы.

Алгоритм использует понятие индекса, который определяется условиями

$\nu = 1$, если $g_1(x) > 0$;

$1 < \nu \leq m$, когда $g_j(x) \leq 0$, но $g_\nu(x) > 0, 1 \leq j \leq \nu - 1$;

$\nu = m + 1$, если $g_j(x) \leq 0, 1 \leq j \leq m$.

Каждая итерация предлагаемого алгоритма включает определение индекса $1 \leq \nu(x_i) \leq m + 1$ точки $x_i, 1 \leq i \leq k$, равного номеру первого нарушенного ограничения. Если $\nu(x_i) < m + 1$, то точке испытания x_i соответствует значение $z_i = g_\nu(x_i)$. Если $\nu(x_i) = m + 1$, (т.е. все ограничения вида $g_i(x) \leq 0$ выполняются), то точке испытания x_i соответствуют значения критерия $z_i = \varphi(x_i)$.

Граничным точкам присваиваются нулевые индексы, значения функций в них не вычисляются. Первая итерация осуществляется в произвольной внутренней точке отрезка $[a, b]$.

Выбор точки $x^{k+1}, k \geq 1$ любой следующей итерации определяется правилами:

1) точки x^1, \dots, x^k предшествующих итераций перенумеровываются нижними индексами в порядке возрастания координаты, т.е.

$$a = x_0 < x_1 < \dots < x_i < x_k < x_{k+1} = b;$$

2) определяются множества

$$I_0 = \{0, k + 1\},$$

$$I_\nu = \{i : 1 \leq i \leq k, \nu = \nu(x_i)\},$$

содержащие номера всех точек, индекс которых равен ν ;
множества

$$S_\nu = \{I_0 \cup \dots \cup I_{\nu-1}\}, 1 \leq \nu \leq m + 1,$$

содержащие номера всех точек, индексы которых меньше ν ;
множества

$$T_\nu = \{I_{\nu+1} \cup \dots \cup I_{m+1}\}, 1 \leq \nu \leq m + 1,$$

содержащие номера всех точек, индексы которых больше ν ;

3) вычисляются максимальные абсолютные значения относительных первых разностей,

$$\mu_\nu = \max\{|z_i - z_p| / (x_i - x_p), i, p \in I_\nu, i > p\}, 1 \leq \nu \leq m + 1;$$

причём в случаях, когда $\text{card } I_\nu < 2$, $1 \leq \nu \leq m+1$ или когда μ_ν оказываются равными нулю, то принимается, что $\mu_\nu = 1$;

4) для всех непустых множеств I_ν , $1 \leq \nu \leq m$ определяются величины

$$z_\nu^* = \begin{cases} 0, & T_\nu \neq 0; \\ \min\{z_i : i \in I_\nu\}, & T_\nu = 0; \end{cases}$$

5) для каждого интервала (x_{i-1}, x_i) , $1 \leq i \leq k+1$ вычисляется характеристика $R(i)$, ($r > 1$ -параметры метода)

$$R(i) = \begin{cases} (x_i - x_{i-1}) + \frac{(z_i - z_{i-1})^2}{(\mu_\nu)^2 (x_i - x_{i-1})^2} - \frac{2(z_i + z_{i-1} - 2z_\nu^*)}{r\mu_\nu}, & v(x_{i-1}) = v(x_i), \\ 2(x_i - x_{i-1}) - 4(z_i - z_\nu^*)/(r\mu_\nu), & v(x_{i-1}) < v(x_i), \\ 2(x_i - x_{i-1}) - 4(z_{i-1} - z_\nu^*)/(r\mu_\nu), & v(x_i) < v(x_{i-1}), \end{cases}$$

6) определить интервал (x_{t-1}, x_t) , имеющий максимальную характеристику, т.е.

$$R(t) = \max\{R(i) : 1 \leq i \leq k+1\};$$

8) очередная итерация осуществляется в точке

$$x^{k+1} = \begin{cases} (x_t + x_{t-1})/2, & v(x_{t-1}) \neq v(x_t); \\ (x_t + x_{t-1})/2 - (z_t - z_{t-1})/(2r\mu_\nu), & v(x_{t-1}) = v(x_t); \end{cases}$$

Представленный индексный алгоритм также является характеристическим, и для него справедлива следующая

Теорема 2.6. Если выполняются условия

1) области Q_i , $1 \leq i \leq M$ являются объединениями конечного числа отрезков положительной длины;

2) функции $g_i(x)$, $1 \leq i \leq M$ допускают удовлетворяющие условиям Липшица продолжения на весь интервал $[a, b]$;

3) для величин μ_ν , начиная с некоторого шага, справедливы неравенства

$$r \mu_\nu > 4 L_\nu, \quad 1 \leq \nu \leq M,$$

то множество предельных точек последовательности $\{x^k\}$, порождаемой алгоритмом, совпадает с множеством решений задачи (2.60), причем индекс каждой предельной точки равен M .

Доказательство дано в работе [37].

Алгоритм можно дополнить условием остановки (по заданной точности $\varepsilon > 0$), прекращающим итерации при выполнении неравенства

$$x_t - x_{t-1} \leq \varepsilon.$$

Краткий обзор главы

В данной главе рассматриваются задачи поиска экстремальных значений функций одной переменной (задачи одномерной оптимизации) и теоретические основы численных методов анализа таких задач. Вводятся различные постановки задач оптимизации, формируется абстрактная модель метода поиска экстремума и приводится общая вычислительная схема алгоритма оптимизации.

Обсуждаются общие свойства методов анализа экстремальных задач (сходимость, оценка погрешности по конечному числу испытаний, значение априорной информации). Особое внимание уделено обсуждению вопросов эффективности методов оптимизации. Сопоставлены различные подходы к формулировке принципов оптимальности алгоритмов поиска экстремума и приведены различные понятия оптимальности алгоритма (оптимальный, ε -оптимальный, последовательно оптимальный, одношагово-оптимальный, асимптотически оптимальный).

Рассмотрен широкий класс методов одномерной оптимизации – класс характеристических алгоритмов, объединяющий в своем составе многие известные методы поиска экстремума. В частности, показано, что любой одношагово-оптимальный алгоритм оптимизации является характеристическим.

Построена общая теория сходимости методов данного класса. Установлены условия двусторонней сходимости к предельным точкам, что позволило обосновать условие останова для характеристических алгоритмов. Получены условия всюду плотной сходимости. Сформулированы требования, при которых характер сходимости алгоритмов рассматриваемого класса имеет иную природу по сравнению со всюду плотной сходимостью, а именно, обеспечивает сходимость только к точкам, обладающим свойством локальной минимальности. Наконец, установлены достаточные условия сходимости ко всем глобально-оптимальным решениям.

В завершение главы рассматривается специальный алгоритм решения задач с ограничениями (индексный метод) и устанавливаются теоретические условия его сходимости к искомому глобально-оптимальному решению.

Глава 3. Фундаментальные способы редукции размерности. Многошаговая схема

3.1. Принципы редуцирования сложности в задачах принятия решений.

Рассмотрим *конечномерную задачу оптимизации* (2.6) в формулировке, которую принято называть *задачей нелинейного программирования*:

$$f(y) \rightarrow \inf, y \in Q \subseteq R^N \quad (3.1)$$

$$Q = \{y \in D : g_j(y) \leq g_j^+, 1 \leq j \leq m\} \quad (3.2)$$

$$D = \{y \in R^N : y_i \in [a_i, b_i], 1 \leq i \leq N\} \quad (3.3)$$

т.е. задачу отыскания экстремальных (в смысле постановок А-Д главы 2) значений целевой (минимизируемой) функции $f(y)$ в области Q , задаваемой *координатными* (3.3) и *функциональными* (3.2) *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 (3.1)-(3.3) является формулировкой *общей многомерной задачи нелинейного программирования*.

Если $m = 0$, т.е. функциональные ограничения отсутствуют, будем полагать $Q = D$. Задача (3.1)-(3.3) в этом случае будет называться задачей *безусловной оптимизации*.

Предметом рассмотрения настоящей главы являются многоэкстремальные задачи оптимизации, т.е. задачи, в которых целевая функция $f(y)$ может иметь в допустимой области Q 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 (либо реализовать) оптимальные (в том или ином смысле алгоритмы) удастся в исключительных случаях.

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 одноэкстремальна. Тогда, поместив начальную точку поиска в эту окрестность, можно тем или иным локальным методом найти искомый 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 [1, 2], т.е. случайно в соответствии с некоторым распределением в области поиска (обычно равномерным), либо использовать в качестве таких точек узлы некоторой регулярной сетки [3]. Сравнение способов выбора начальных точек приведено в [4].

При таком способе сходимости 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, множество базовых точек может модифицироваться в процессе поиска – туда могут добавляться новые точки. Также в процессе реализации этой схемы может изменяться способ 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

$$\lim_{\tau \rightarrow \infty} \min_{1 \leq k \leq \tau} f(y^k) = \min_{y \in Q} f(y) \quad (3.4)$$

Другими словами, последовательность испытаний этих методов является всюду плотной в области поиска, т.е. сходится в смысле Определения 2.2 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 плодотворно.

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 [5-9] и редукция размерности на основе кривых, заполняющих пространство (кривых Пеано) [6-8, 10].

Настоящая глава будет посвящена рассмотрению первой из указанных схем – многошаговой схемы редукции размерности, в которой решение задачи (3.1)-(3.3) сводится к решению семейства рекурсивно связанных одномерных подзадач.

3.2. Многошаговая схема редукции размерности

Рассмотрим в качестве исходной задачу (3.1)-(3.3) 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 $\tilde{g}_j(y) = g_j(y) - g_j^+ \leq 0$), 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 Q .

Введем непрерывную функцию, определенную в области D , такую, что

$$\begin{aligned} G(y) &\leq 0, y \in Q \\ G(y) &> 0, y \notin Q \end{aligned} \quad (3.5)$$

В качестве $G(y)$ можно взять, например,

$$G(y) = \max\{g_1(y), \dots, g_m(y)\} \quad (3.6)$$

или

$$G(y) = \max\{0; g_1(y), \dots, g_m(y)\} \quad (3.7)$$

Последняя функция тождественно равна нулю в области Q .

Введем обозначения

$$u_i = (y_1, \dots, y_i), \quad v_i = (y_{i+1}, \dots, y_N), \quad (3.8)$$

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 Q :

$$S_1 = Q, \quad S_{i+1}(u_i) = \{(u_i, v_i) \in Q\}, \quad 1 \leq i \leq N-1, \quad (3.9)$$

и проекции сечений на ось y_{i+1} :

$$\Pi_{i+1}(u_i) = \{y_{i+1} \in R^1 : \exists (y_{i+1}, v_{i+1}) \in S_{i+1}(u_i)\}, \quad (3.10)$$

Положим $G^N(y) \equiv G(y)$ и построим семейство функций

$$G^i(u_i) = \min\{G^{i+1}(u_i, y_{i+1}) : y_{i+1} \in [a_{i+1}, b_{i+1}]\}, \quad 1 \leq i \leq N-1, \quad (3.11)$$

определенных в соответствующих проекциях

$$D_i = \{u_i \in R^i : y_j \in [a_j, b_j], 1 \leq j \leq i\} \quad (3.12)$$

множества D из (3.3) 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 далее определяет существование и непрерывность всех функций семейства (3.11).

Справедлива следующая лемма.

Лемма 3.1.

$$G^i(u_i) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\}, 1 \leq i \leq N-1. \quad (3.13)$$

Доказательство. Вследствие (3.11) доказательство (3.13) состоит в установлении справедливости равенства

$$\min_{y_{i+1} \in [a_{i+1}, b_{i+1}]} \dots \min_{y_N \in [a_N, b_N]} G(u_i, v_i) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\} \quad (3.14)$$

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 (3.14) 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 (3.11) 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

$$Q_i = \{u_i \in R^i : \exists(u_i, v_i) \in Q\}, 1 \leq i \leq N, \quad (3.15)$$

множества Q на координатные оси y_1, \dots, y_i .

Лемма 3.2. Представление (3.15) эквивалентно соотношению

$$Q_i = \{u_i \in R^i : G^i(u_i) \leq 0\} \tag{3.16}$$

Доказательство. Пусть выполнено (3.15), т.е. для некоторого u_i существует v_i^* такой, что $(u_i, v_i^*) \in Q$. Но тогда $G(u_i, v_i^*) \leq 0$, т.е. вследствие Леммы 3.1

$$G^i(u_i) = \min\{G(u_i, v_i) : y_j \in [a_j, b_j], i+1 \leq j \leq N\} \leq G(u_i, v_i^*) \leq 0,$$

и (3.16) справедливо.

Пусть теперь, наоборот, для некоторого $u_i \in Q_i$ выполняется $G^i(u_i) \leq 0$. Но согласно Лемме 3.1 существует вектор v_i^* такой, что $G^i(u_i) = G(u_i, v_i^*)$, т.е. $(u_i, v_i^*) \in Q$. Лемма доказана.

Лемма 3.3. Определение проекции (3.10) эквивалентно представлению

$$\Pi_{i+1}(u_i) = \{y_{i+1} \in [a_{i+1}, b_{i+1}] : G^{i+1}(u_i, y_{i+1}) \leq 0\}, \tag{3.17}$$

Доказательство. Прежде всего заметим, что поскольку $Q \subseteq D$, 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 (3.10).

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 $(u_i, y_{i+1}, v_{i+1}^*) \in Q$, 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

$$f^i(u_i) = \min\{f^{i+1}(u_i, y_{i+1}) : y_{i+1} \in \Pi_{i+1}(u_i)\}, 1 \leq i \leq N-1, \tag{3.18}$$

определенных на соответствующих проекциях Q_i . Тогда имеет место основное соотношение

$$\min_{y \in Q} f(y) = \min_{y_1 \in \Pi_1} \min_{y_2 \in \Pi_2(u_1)} \dots \min_{y_N \in \Pi_N(u_{N-1})} f(y) \tag{3.19}$$

Как следует из (3.19), для решения задачи (3.1) – (3.3) достаточно решить одномерную задачу

$$f^1(y_1) \rightarrow \min, y_1 \in \Pi_1 \subseteq R^1 \tag{3.2 0}$$

$$\Pi_1 = \{y_1 \in [a_1, b_1] : G^1(y_1) \leq 0\} \tag{3.2 1}$$

При этом каждое вычисление функции $f^1(y_1)$ в некоторой фиксированной точке $y_1 \in \Pi_1$ представляет собой согласно (3.18) решение одномерной задачи

$$f^2(y_1, y_2) \rightarrow \min, y_2 \in \Pi_2(y_1) \subseteq R^1 \tag{3.2 2}$$

$$\Pi_2(y_1) = \{y_2 \in [a_2, b_2] : G^2(y_1, y_2) \leq 0\} \tag{3.2 3}$$

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

$$f^3(u_2, y_3) \rightarrow \min, y_3 \in \Pi_3(u_2) \tag{3.2 4}$$

и т.д. вплоть до решения задачи

$$f^N(u_{N-1}, y_N) = f(u_{N-1}, y_N) \rightarrow \min, y_N \in \Pi_N(u_{N-1}) \tag{3.2 5}$$

при фиксированном u_{N-1} .

Окончательно решение задачи (3.1) – (3.3) сводится к решению семейства "вложенных" одномерных подзадач

$$f^i(u_{i-1}, y_i) \rightarrow \min, y_i \in \Pi_i(u_{i-1}), \tag{3.2 6}$$

где фиксированный вектор $u_{i-1} \in Q_{i-1}$.

Решение исходной многомерной задачи (3.1) – (3.3) через решение системы взаимосвязанных одномерных подзадач (3.26) называется *многошаговой схемой редукции размерности*.

Заметим, что если в задачах (3.26) по некоторым координатам y_i искать не минимум, а максимум, то становится возможным вычисление сложных минимаксных (максиминных) выражений вида

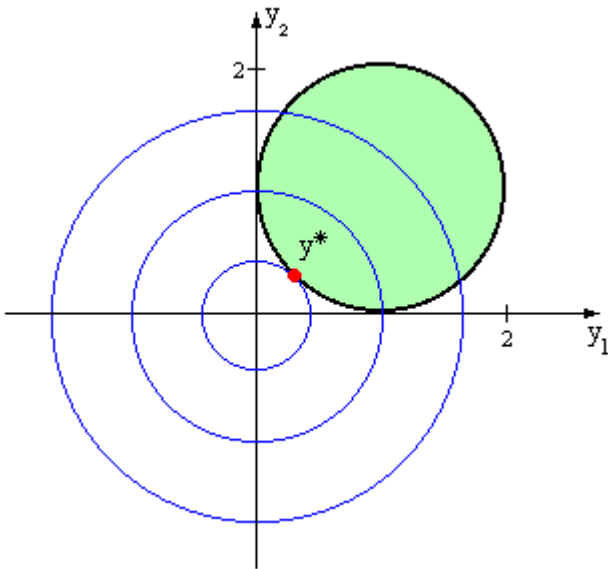
$$\begin{matrix} \text{extr} & f(y) = & \text{extr} & & \text{extr} & \dots & \text{extr} & f(y), \\ y \in Q & & y_1 \in \Pi_1 & y_2 \in \Pi_2(u_1) & y_N \in \Pi_N(u_{N-1}) & & & \end{matrix}$$

где операция *extr* означает вычисление глобального минимума или максимума. Проиллюстрируем общие результаты конкретным несложным примером.

Пример 3.1. Рассмотрим двумерную задачу (3.1) – (3.3), в которой

$$f(y) = y_1^2 + y_2^2, \quad Q = \{y \in D : (y_1 - 1)^2 + (y_2 - 1)^2 - 1 \leq 0\},$$

$$D = \{y \in R^2 : 0 \leq y_1, y_2 \leq 2\}. \tag{3.27}$$



На рисунке серым цветом отмечена допустимая область, а также показаны концентрические 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 (3.17) 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

$$\Pi_2(y_1) = [1 - \sqrt{1 - (y_1 - 1)^2}, 1 + \sqrt{1 - (y_1 - 1)^2}] \tag{3.28}$$

Функция $f(y) = y_1^2 + y_2^2$, будучи возрастающей по y_2 на отрезке $[0, 2]$ и, следовательно, в области (3.28), 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 Q . 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 (3.27) 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 точка оптимума помечена темным кружком.

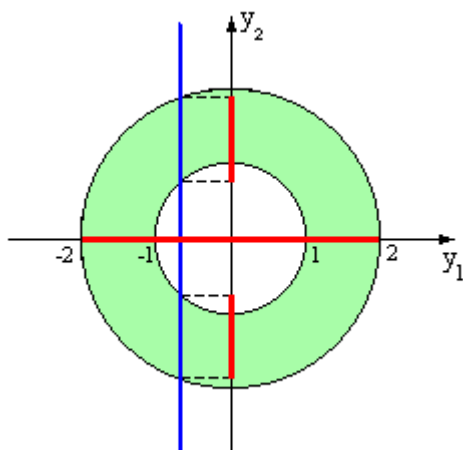
В рассмотренном примере мы построили границы областей одномерного поиска аналитически, установив области неположительности соответствующих функций $G^i(u_i)$. Вместе с тем можно указать более наглядный "геометрический" способ построения проекций $\Pi_{i+1}(u_i)$. Собственно, этот способ вытекает из определений (3.9) и (3.10) и состоит в том, что необходимо построить сечения области Q 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 (3.19) аналогично процедуре нахождения многомерного интеграла от функции $f(y)$ в области Q 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.

Пример 3.2. Пусть область оптимизации (или интегрирования) задается как

$$Q = \{y \in R^2 : -4 \leq y_1, y_2 \leq 4, y_1^2 + y_2^2 \leq 4, y_1^2 + y_2^2 \geq 1\} \quad (3.29)$$



Серый цвет, как и ранее, помечает допустимую область. Левая 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, а также проекция Π_1 на рисунке отмечены красным цветом.

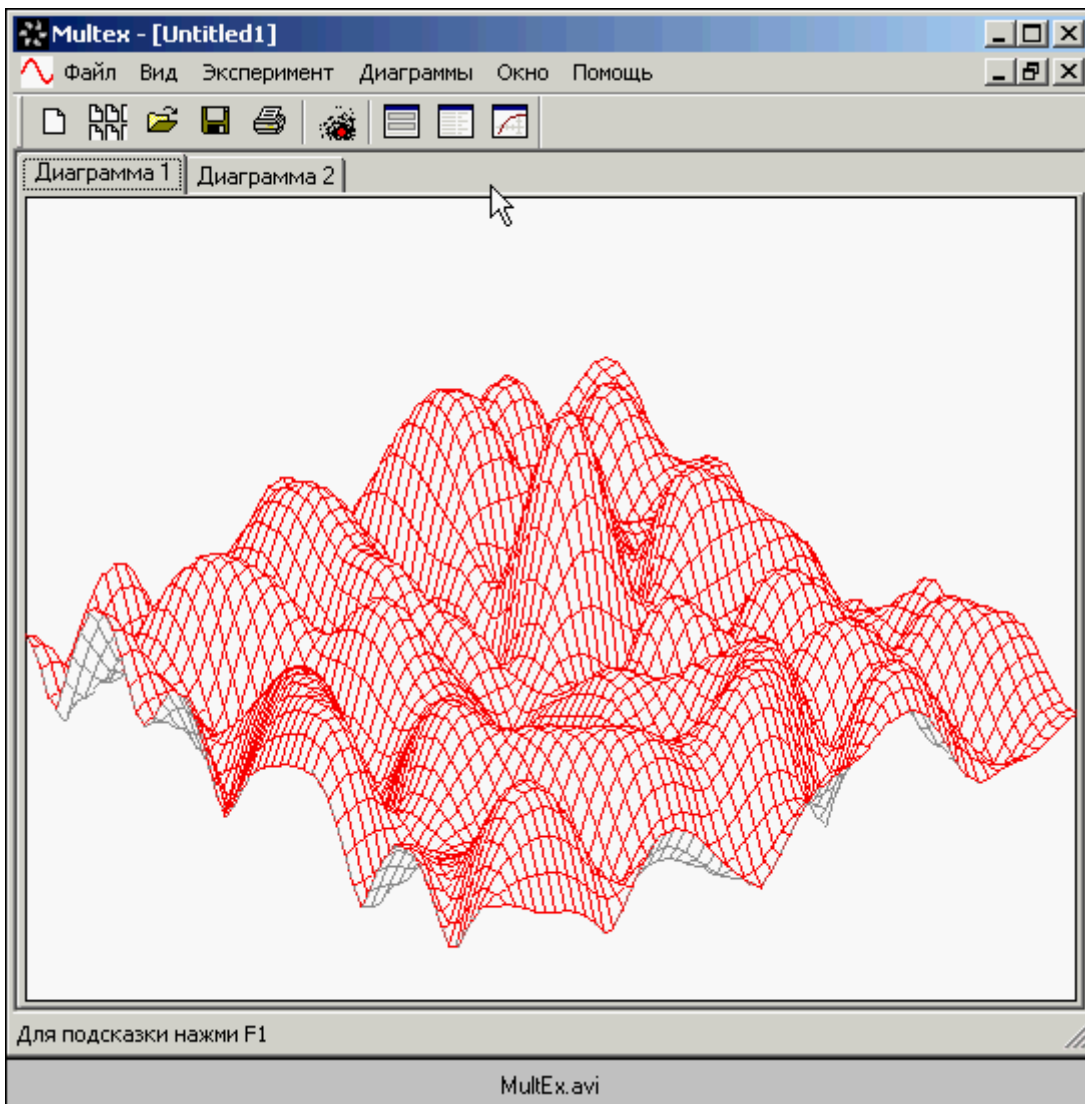
Подобные геометрические соображения позволяют построить требуемые проекции задачи как

$$\Pi_1 = [-2, 2],$$

$$\Pi_2(y_1) = \begin{cases} [-\sqrt{4 - y_1^2}, \sqrt{4 - y_1^2}], y_1 \in [-2, -1] \cup [1, 2]; \\ [-\sqrt{4 - y_1^2}, -\sqrt{1 - y_1^2}] \cup [\sqrt{1 - y_1^2}, \sqrt{4 - y_1^2}], y_1 \in [-1, 1]. \end{cases}$$

Пример 3.3.

Данный пример демонстрирует процесс максимизации двумерной многоэкстремальной функции в области, задаваемой сложными ограничениями. Сначала демонстрируется поверхность целевой функции (без учета ограничений), а затем линии уровня данной функции с отображением допустимой области (сиреневый цвет) и недопустимых участков (помечены розовым цветом). На диаграмме линий уровня показывается точка и величина глобального максимума задачи и запускается процесс решения по многошаговой схеме с использованием метода Стронгина при решении одномерных подзадач. Координаты испытаний отмечаются крестиками.



Контрольные вопросы и упражнения:

1. Постройте проекции Π_1 и $\Pi_2(y_1)$ для области $Q = \tilde{Q} \cup \hat{Q}$, 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 (3.29).

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 (3.18). 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 Q . 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 Q , определяемые ограничениями (3.2.), (3.3).

3.3. Свойства одномерных подзадач многошаговой схемы

3.3.1. Структура допустимых областей одномерного поиска

Для анализа структуры областей $\Pi_i(u_{i-1})$ используем результаты Леммы 3.3, которая установила эквивалентность определения (3.10) и представления (3.17). Дело в том, что (3.17) 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 (3.12), 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 (3.26) является задачей вида

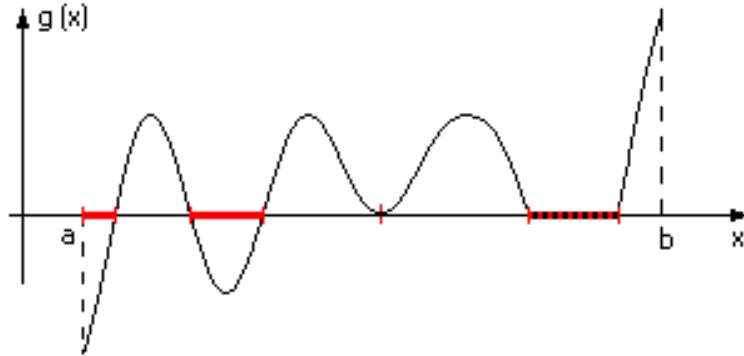
$$\begin{aligned} \varphi(x) \rightarrow \min, x \in \bar{Q} \subset R^1, \\ \bar{Q} = \{x \in [a, b] : g(x) \leq 0\}, \end{aligned} \quad \begin{matrix} (3.3) \\ 0) \end{matrix}$$

причем функция $g(x)$ непрерывна.

Непрерывность ограничения $g(x)$ позволяет утверждать, что допустимая область \bar{Q} может быть записана в виде системы отрезков

$$\bar{Q} = \bigcup_{j=1}^q [a^j, b^j] \tag{3.3 1)}$$

на каждом из которых функция неположительна. Для примера рассмотрим рисунок, который отражает возможные случаи поведения непрерывной функции, порождающие области неположительности в виде отрезков (помечены утолщениями на оси x), включая касание оси x в точке. Самый правый отрезок, отмеченный пунктиром,



соответствует ситуации, когда функция на данном отрезке равна нулю во всех его точках.

В системе (3.31) 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 (3.26) есть множество вида (3.31), т.е.

$$\Pi_i(u_{i-1}) = \bigcup_{j=1}^{q_i} [a_i^j, b_i^j], \tag{3.3 2)}$$

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.

$$q_i = q_i(u_{i-1}), \quad a_i^j = a_i^j(u_{i-1}), \quad b_i^j = b_i^j(u_{i-1}) \tag{3.3 3)}$$

Если область Q такова, что для всех $1 \leq i \leq N$ удается указать явные (аналитические) выражения для величин q_i, a_i^j, b_i^j как функций $u_{i-1} \in Q_{i-1}$, тогда область Q называется *областью с вычислимой границей*. Образцы таких областей приведены в примерах 3.1 и 3.2. Для построения данных областей необходимо уметь аналитически находить все корни функций $G^i(u_{i-1}, y_i)$ по соответствующим переменным 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 (3.30) Для отыскания всех корней непрерывной функции $g(x)$ можно предложить численное решение эквивалентной задачи оптимизации

$$|g(x)| \rightarrow \min, x \in [a, b], \quad (3.34)$$

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 (3.34), рассмотрен в параграфе 2.5 предыдущей главы.

Практически важным частным случаем задачи (3.1) – (3.3) является случай $Q = D$, когда функциональные ограничения (3.2) отсутствуют. В данной ситуации $G(y) \equiv 0$ в области D , а из (3.11) следует, что функции $G^i(u_i) \equiv 0$, $u_i \in D_i$. Тогда согласно (3.17)

$$\Pi_i(u_{i-1}) = [a_i, b_i], \quad (3.35)$$

где a_i, b_i - константы.

Другим важным частным случаем является случай *выпуклых* ограничений.

Определение 3.1. Функция $g(y)$ называется выпуклой (вниз) в выпуклой области Q , если для любых $y', y'' \in Q$ и для любых $\alpha \in [0,1]$ выполняется

$$g(\alpha y' + (1 - \alpha)y'') \leq \alpha g(y') + (1 - \alpha)g(y'') \quad (3.36)$$

Теорема 3.1. Если в задаче (3.1) – (3.3) ограничения $g_j(y)$, $1 \leq j \leq m$, выпуклы, функция $G(y)$ выбирается согласно (3.6) или (3.7), тогда проекции $\Pi_{i+1}(u_i)$ из (3.10) либо пусты, либо имеют вид

$$\Pi_{i+1}(u_i) = [a_{i+1}^1(u_i), b_{i+1}^1(u_i)], \quad (3.37)$$

т.е. область одномерного поиска в (3.26) состоит из одного отрезка.

Доказательство.

1). Покажем вначале, что область Q является выпуклой. Возьмем произвольные $y', y'' \in Q$ 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, $\alpha y' + (1-\alpha)y'' \in Q$ при любом $\alpha \in [0,1]$, что и доказывает выпуклость Q .

2). Теперь убедимся, что проекции Q_i из (3.15) также выпуклы.

Возьмем два произвольных вектора $u'_i, u''_i \in Q_i$. Тогда существуют вектора v'_i, v''_i такие, что $y' = (u'_i, v'_i), y'' = (u''_i, v''_i) \in Q$. Но тогда для любого $\alpha \in [0,1]$ вектор $\alpha y' + (1-\alpha)y'' = (\alpha u'_i + (1-\alpha)u''_i, \alpha v'_i + (1-\alpha)v''_i) \in Q$, т.е. вектор $\alpha u'_i + (1-\alpha)u''_i \in Q_i$.

3). Покажем, что функция (3.6) выпукла в D .

Для любой точки $\alpha y + (1-\alpha)z$, $\alpha \in [0,1]$, $y, z \in D$, существует номер s такой, что

$$G(\alpha y + (1-\alpha)z) = g_s(\alpha y + (1-\alpha)z) \leq \alpha g_s(y) + (1-\alpha)g_s(z) \\ \leq \alpha \max\{g_1(y), \dots, g_m(y)\} + (1-\alpha) \max\{g_1(z), \dots, g_m(z)\} = \alpha G(y) + (1-\alpha)G(z)$$

4). Докажем выпуклость функций $G^i(u_i)$ из (3.11) 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

$$G^i(u'_i) = G^{i+1}(u'_{i+1}), \quad u'_{i+1} = (u'_i, y'_{i+1}), \\ G^i(u''_i) = G^{i+1}(u''_{i+1}), \quad u''_{i+1} = (u''_i, y''_{i+1}),$$

Выберем произвольное $\alpha \in [0,1]$ и обозначим $y_{i+1}^* = \alpha y'_{i+1} + (1-\alpha)y''_{i+1}$ и $u_i^\alpha = \alpha u'_i + (1-\alpha)u''_i$. Тогда существует $y_{i+1}^\alpha \in [a_{i+1}, b_{i+1}]$ такой, что

$$G^i(u_i^\alpha) = G^{i+1}(u_i^\alpha, y_{i+1}^\alpha) \leq G^{i+1}(u_i^\alpha, y_{i+1}^*) \leq \alpha G^{i+1}(u'_i, y'_{i+1}) + (1-\alpha)G^{i+1}(u''_i, y''_{i+1}) = \\ = \alpha G^i(u'_i) + (1-\alpha)G^i(u''_i)$$

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 [11], 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 (3.37) не только в случае выпуклых ограничений, но и в более общем случае монотонно унимодальных ограничений [12], которые могут порождать невыпуклые множества Q .

Определение 3.2. 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 3.3. 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 возрастая на втором.

Определение 3.4. Ограничения (3.2) 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 (3.2) невыпуклую допустимую область Q .

Теорема 3.2. Пусть ограничения $g_j(y)$ из (3.2) монотонно унимодальны в совокупности в области D и функция $G(y)$ строится согласно (3.6) или (3.7). 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 (3.10) имеет вид (3.37).

Доказательство теоремы может быть найдено в [12].

3.3.2. Свойства целевых функций в одномерных подзадачах

Целевой функцией в подзадаче (3.26) является функция $f^i(u_{i-1}, y_i)$ при фиксированном u_{i-1} , и для решения подзадач (3.26) определяющим является характер зависимости функции f^i от переменной y_i .

Рассмотрим класс задач, в которых функция $f(y)$ является сепарабельной, т.е.

$$f(y) = \sum_{i=1}^N f_i(y_i) \quad (3.38)$$

а функциональные ограничения $g_j(y)$ отсутствуют, т.е. $Q = D$.

Тогда, как следует из (3.19), (3.38),

$$\min_{y \in Q} f(y) = \sum_{i=1}^N \min_{y_i \in [a_i, b_i]} f_i(y_i),$$

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 Q условию Липшица с константой $L > 0$, т.е. для любых $y', y'' \in Q$

$$|f(y') - f(y'')| \leq L \|y' - y''\|. \quad (3.39)$$

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 (3.1) – (3.3) целевая функция $f(y)$ удовлетворяет условию Липшица (3.39), а область Q имеет вид

$$Q = \{y \in R^2 : y_1^2 + y_2^2 - 1 \leq 0\}$$

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 (3.38) уже не подчиняется. В [6] 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 [6] устанавливает достаточные условия липшицевости функций $f^i(u_i)$.

Теорема 3.3. Пусть функция $f(y)$ является липшицевой с константой L в выпуклой области Q из (3.2) и граничные пары (3.37) являются кусочно-линейными функциями вида

$$a_{i+1}(u_i) = \max_{1 \leq \nu \leq p_i} \{\alpha_i^\nu u_i + A_i^\nu\} \quad (3.4 \ 0)$$

$$b_{i+1}(u_i) = \max_{1 \leq \nu \leq r_i} \{\beta_i^\nu u_i + B_i^\nu\}, \quad (3.4 \ 1)$$

где $\alpha_i^\nu u_i$, $\beta_i^\nu u_i$ - есть скалярные произведения векторов из R^i и A_i^ν, B_i^ν - константы. Тогда функции $f^i(u_i)$, $u_i \in Q_i$, являются липшицевыми с константами L_i , $1 \leq i \leq N$, где

$$L_N = L, \quad L_i = L \prod_{j=i}^{N-1} (1 + \lambda_j), \quad 1 \leq i \leq N-1,$$

$$\lambda_j = \max \left\{ \max_{1 \leq \nu \leq p_i} \|\alpha_j^\nu\|, \max_{1 \leq \nu \leq r_i} \|\beta_j^\nu\| \right\}$$

Доказательство теоремы приведено в [6].

Комментарии к теореме.

1. Представление граничных пар в виде (3.40), (3.41) имеет место, если область является выпуклым многогранником.

2. Если область Q является гиперпараллелепипедом, то все вектора $\alpha_i^\nu, \beta_i^\nu$ 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 Q .

В самом простом случае, когда множество Q 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.

Краткий обзор главы

Настоящая глава посвящена проблематике исследования многомерных многоэкстремальных задач оптимизации при наличии ограничений. Отмечается высокая вычислительная сложность таких задач и дается обзор основных подходов к их эффективному решению, основанных на принципах редукции сложности. Одним из базовых направлений в данном контексте является применение схем редукции

размерности, в частности, многошаговой схемы редукции, рассмотрение которой является основным содержанием главы.

С этой целью вводятся базовые понятия, которые используются для формулировки основного соотношения многошаговой схемы, и рассматривается структура редуцирования многомерной задачи к системе взаимосвязанных одномерных подзадач.

Обсуждаются способы построения и свойства одномерных подзадач, порождаемых многошаговой схемой редукции. Устанавливается структура областей одномерного поиска как системы отрезков. Формулируются условия (отсутствие ограничений, выпуклость, монотонная унимодальность), при которых области поиска решения в одномерных подзадачах представляются в виде единственного отрезка.

Изучается вопрос о связи характеристик исходной задачи со свойствами целевых функций одномерных подзадач. Дается анализ влияния ограничений на сложность одномерных подзадач для случая липшицевости многомерной оптимизируемой функции.

Глава 4. Модели и методы поиска локально-оптимальных решений

4.1. Постановка задачи поиска локально-оптимальных решений

В предыдущих главах уже отмечалось, что в задачах математического программирования

$$f(y) \rightarrow \min, y \in Q \subseteq R^N, \quad (4.1)$$

$$Q = \{y \in D : g_j(y) \leq g_j^+, j = 1, \dots, m\}, \quad (4.2)$$

$$D = \{y \in R^N : a_i \leq y_i \leq b_i, i = 1, \dots, N\}, \quad (4.3)$$

возникающих в рамках общей модели рационального выбора, используются два понятия решения, глобальное и локальное.

Всякий локальный минимум $\hat{y} \in Q$ в задаче математического программирования определяет ее *локальное (локально-оптимальное) решение*.

Глобальный минимум $y^* \in Q$, являясь одним из локальных, характеризует наилучшее (по значению целевой функции) из локально-оптимальных решений задачи.

Напомним, что точка $\hat{y} \in Q$ называется точкой *локального минимума* функции $f(y)$ на множестве Q , если существует такое число $\varepsilon > 0$, что для всех $y \in Q$ таких, что $|y - \hat{y}| < \varepsilon$ выполняется $f(\hat{y}) \leq f(y)$.

Задача поиска локально-оптимального решения возникает, как правило, тогда, когда известна приближенная оценка y^0 глобально-оптимального решения, найденная с неудовлетворительной точностью. В этом случае достаточно найти с высокой точностью локально-оптимальное решение, соответствующее начальной точке поиска y^0 . Если эта точка была выбрана правильно, то найденный локальный минимум $\hat{y}(y^0)$, зависящий от y^0 , 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, что и 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 третьей главе). В задачах высокой размерности практически единственным средством их решения остаются методы локальной оптимизации,

совмещенные, как это было описано в разделе 3.1, с процедурами предварительного отбора начальных точек, используемых при локальной оптимизации.

Проблема локальной оптимизации заключается в том, чтобы в задаче математического программирования (4.1) по заданной начальной точке $y^0 \in Q$ определить локальный минимум $\hat{y}(y^0)$ со значением целевой функции, не превосходящим $f(y^0)$.

Прежде чем рассматривать методы локальной оптимизации для задач с ограничениями в постановке (4.1)–(4.3), следует изучить необходимый материал по методам локальной оптимизации в более простых задачах, не содержащих ограничений. Этому вопросу помещены второй и третий разделы четвертой главы. В следующих за ними разделах приведены модификации методов локальной оптимизации применительно к задачам с ограничениями.

Итак, вначале рассмотрим задачи без ограничений

$$f(y) \rightarrow \min, y \in R^N. \quad (4.4)$$

В них кроме целевой функции $f(y)$, определенной для $y \in R^N$, будем также задавать начальную точку y^0 . Предполагается, что минимум в задаче существует. Требуется определить локальный минимум функции $f(y)$, соответствующий начальной точке y^0 .

В дальнейшем, для удобства изложения, будем дополнительно предполагать, что задача (4.4) является *одноэкстремальной*, т.е. имеет единственный локальный минимум, который в этом случае $\hat{y}(y^0)$ является точкой y^* ее глобального минимума.



Замечание. Для одноэкстремальных задач при дальнейшем изложении в обозначениях не будет проводиться различий между локальным и глобальным минимумами. Всегда будет использоваться обозначение y^* .

4.2. Общие принципы построения методов локальной оптимизации

4.2.1. Структура методов поиска локального минимума функций

К настоящему времени разработано огромное количество разнообразных методов локальной оптимизации. Большинство из них реализует идею *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

$$y^{k+1} = y^k + x^k d^k, \quad (4.5)$$

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 (гессиан) $\Gamma_k = \Gamma(y^k)$. 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 .

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 $f^{k+1} \leq f^k$, для обеспечения сходимости к решению задачи недостаточно.

Останов вычислений в методах локальной оптимизации, применяемых в задачах без ограничений (4.4) с непрерывно дифференцируемыми целевыми функциями, происходит при выполнении условия достаточной малости нормы градиента

$$\|\nabla f(y^k)\| \leq \varepsilon. \quad (4.6)$$

Нужно отметить, что выполнение условия (4.6), 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.

4.2.2. 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 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 будем называть *моделью задачи*. В большинстве случаев, принятая модель задачи (4.4) может быть описана в терминах принадлежности целевой функции $f(y)$ некоторому классу функций. Будем записывать это следующим образом: $f \in \Phi$.

Принятая модель задачи существенно влияет на интерпретацию результатов проведенных испытаний. Поясним это на нескольких примерах, напомнив предварительно некоторые известные определения и факты, касающиеся выпуклых функций.

Функция $f(y)$, определенная в выпуклой области Q , называется *выпуклой (вниз)*, если для любых двух точек y^1 и y^2 из Q и $\forall \alpha \in [0, 1]$ выполняется неравенство

$$f(\alpha \cdot y^1 + (1 - \alpha)y^2) \leq \alpha \cdot f(y^1) + (1 - \alpha)f(y^2). \quad (4.7)$$

Если неравенство (4.7) является строгим $\forall \alpha \in (0, 1)$, то функция — *строго выпукла*.

Непрерывно дифференцируемая функция $f(y)$, определенная в выпуклой области Q , выпукла тогда и только тогда, когда для любых двух точек y^1 и y^2 из Q выполняется неравенство

$$f(y^2) \geq f(y^1) + (\nabla f(y^1), y^2 - y^1). \quad (4.8)$$

Оно означает, что свойство выпуклости функции равносильно тому, что значения любой ее линейной аппроксимации не превосходят значений самой функции.

ПРИМЕР 1. Пусть функция f принадлежит классу непрерывно дифференцируемых выпуклых (вниз) функций, а ее испытания включают измерение градиента. Тогда по результату испытания в некоторой точке y^k можно построить оценку $\bar{Q}(f^k)$ множества точек $Q(f^k)$, в которых функция $f(y)$ принимает значения, не превосходящие величины $f^k = f(y^k)$. С учетом (4.8), $f^k + (\nabla f^k, y - y^k) \leq f(y)$. Поэтому оценка будет иметь следующий вид

$$Q(f^k) = \{y \in Q : f(y) \leq f^k\} \subseteq \bar{Q}(f^k) = Q \cap \{y \in R^N : f^k + (\nabla f^k, y - y^k) \leq f^k\} = \\ = Q \cap \{y \in R^N : (\nabla f^k, y - y^k) \leq 0\}.$$

Таким образом, проведение каждого испытания позволяет отсеять часть допустимой области, не содержащую решения.

ПРИМЕР 2. Пусть функция f принадлежит классу функций, липшицевых в области Q с константой L^f (функции этого класса рассматривались в главах 2 и 3), а испытание включает только вычисление значения функции, т.е. $I^k = f^k$. Пусть, кроме того, уже известны результаты первых $(k-1)$ -го испытания с результатами вычислений функции f^1, \dots, f^{k-1} . Допустим также, что ранее проведенные испытания позволили сократить исходную область поиска Q до новой области Q^{k-1} , имеющей меньшую меру по сравнению с Q .

Тогда после проведения k -го испытания можно применить следующее правило уточнения оценки множества, содержащего решение задачи

$$Q^k = Q^{k-1} \cap \{y : f^k - L^f \|y - y^k\| \leq \min\{f^j : j = 1, \dots, k\}\}.$$

ПРИМЕР 3. Пусть о функции f известно только то, что она непрерывно дифференцируема, т.е. принадлежит классу $\Phi^C = C(D)$, а испытание в точке y^k состоит в вычислении значения функции f^k и градиента ∇f^k (также, как и в примере 1). В этом случае 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 целесообразно принять направление *антиградиента* $d^k = -\nabla f^k$, определяющего направление скорейшего локального убывания функции.



Замечание. Классы функций, наиболее характерные для задач локальной оптимизации, обычно таковы, что не позволяют сокращать по результатам конечного числа испытаний область возможного положения решения. Именно по этой причине методы локальной оптимизации имеют структуру (4.5), существенно отличающуюся от структуры методов многоэкстремальной оптимизации, рассмотренной в разделах 2.4 и 3.2. Локальный поиск почти всегда основан на выборе направлений существенного локального убывания функции и смещениях вдоль них.

4.2.3. Классификация методов локального поиска

Обычно используют классификацию методов в зависимости от той локальной информации, которую метод получает при выполнении основных испытаний.

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). 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 разделе 4.2.1. Почти все остальные методы соответствуют структуре (4.5), и, следовательно, требуют для своей реализации разработки специальных вычислительных процедур, позволяющих определять в (4.5) коэффициенты одномерных смещений x^k вдоль выбираемых направлений d^k .

4.2.4. Эффективные стратегии поиска вдоль направлений. Регуляризованные алгоритмы одномерного поиска

Вычислительная схема методов локального поиска (4.5) 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 нескольких критериев. Наиболее распространенными при построении методов являются следующие [1], [2]: 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

$$x^k \in \Pi_1(\eta), \quad 0 \leq \eta < 1 \quad (4.9)$$

$$\Pi_1(\eta) = \{x \geq 0 : |(\nabla f(y^k + x \cdot d^k), d^k)| \leq \eta \cdot (-\nabla f(y^k), d^k)\}. \quad (4.10)$$

В основе второго критерия (существенности убывания функции) лежит требование

$$x^k \in \Pi_2(\mu), \quad 0 < \mu < 1 \quad (4.11)$$

$$\Pi_2(\mu) = \{x \geq 0 : f(y^k + x \cdot d^k) \leq f(y^k) + \mu \cdot x \cdot (\nabla f(y^k), d^k)\}. \quad (4.12)$$

Эти два критерия используются совместно, и окончательное условие выбора x^k состоит в одновременном выполнении требований (4.9)-(4.12). При этом для их непротиворечивости на значения параметров μ и η накладывается дополнительное требование вида $\mu < \eta$, т.е.

$$x^k \in \Pi = \Pi_1(\eta) \cap \Pi_2(\mu), \quad 0 < \mu < \eta < 1. \quad (4.13)$$

Напомним, что используемые в формулах (4.10), (4.12) скалярные произведения градиента функции $f(y)$ на вектор направления d определяют производные функции f в точке y в направлении d , а именно $\partial f(y)/\partial d = (\nabla f(y), d)$.

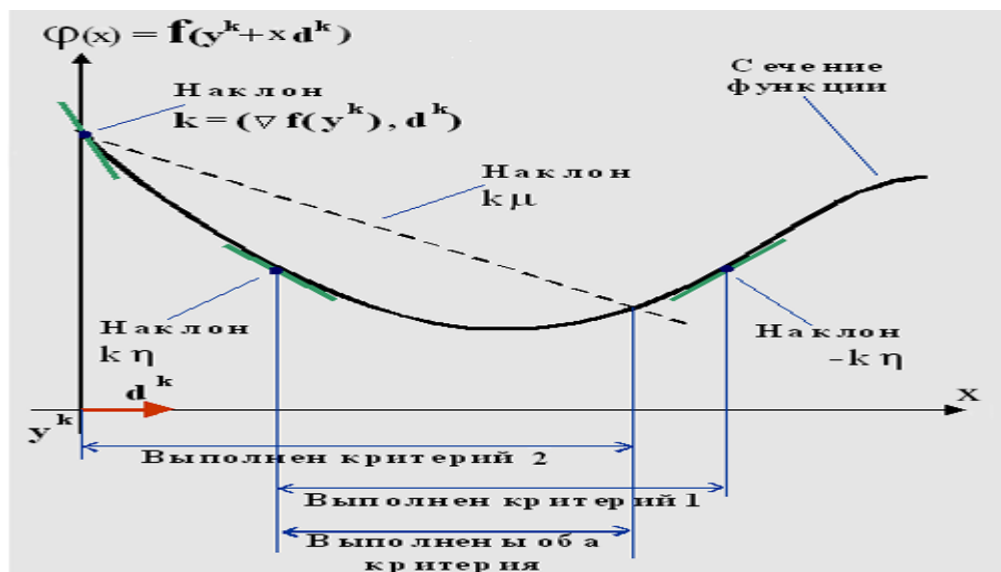


Рис. 4. 1. Критерии выбора коэффициента одномерного шага

На рис.4.1 дана иллюстрация выбора x^k на основе этих двух критериев. Если при выборе коэффициента одномерного шага используется условие (4.10) 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

$$\varphi(x) = f(y^k + x d^k). \quad (4.14)$$

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 ограничений этот промежуток должен иметь пересечение с искомым множеством Π из (4.13). Если же в задаче есть ограничения, то такого пересечения может не существовать, и тогда значение X определяется из условия попадания на границу области D из (4.3).

АЛГОРИТМ ОПРЕДЕЛЕНИЯ ПРОМЕЖУТКА $[0, X]$

ШАГ 0. При решении задачи без ограничений (4.4) выбирается начальное значение $X^0 = \infty$. При решении задачи с ограничениями (4.1)–(4.3), 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 $[0, X]$ осуществляется процедура поиска минимума функции одного переменного $\varphi(x)$ из (4.14), 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 *униmodalной*, что позволяет применить известный *алгоритм золотого сечения*, близкий к ε -оптимальному методу Фибоначчи ([3], [4], а также [5] и [6]). 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 *квазиньютоновскими*. Известно, что при определенных условиях они способны обеспечить скорость сходимости более высокого порядка, чем метод золотого сечения [7]. Если же условия их сходимости будут нарушены, то такие методы могут расходиться.

Для осуществления одномерного поиска гладкой *униmodalной* функции наиболее подходящими являются *регуляризованные алгоритмы*, представляющие комбинацию метода золотого сечения с квазиньютоновским алгоритмом [2].

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 $x_3 = B - (B-A)x$, $\varphi_3 = \varphi(x_3)$. 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. u определяется как точка минимума квадратичной аппроксимации функции $\varphi(x)$, построенной по значениям x_1, φ_1 ; x_2, φ_2 ; x_3, φ_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 (4.13). Если она принадлежит, или же $|B - A| < \sigma X$, 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.

4.3. Классические методы локальной оптимизации

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 [2, 6–10].

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 (4.5), 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.4.2). Более того, это направление не инвариантно по отношению к растяжениям пространства переменных.

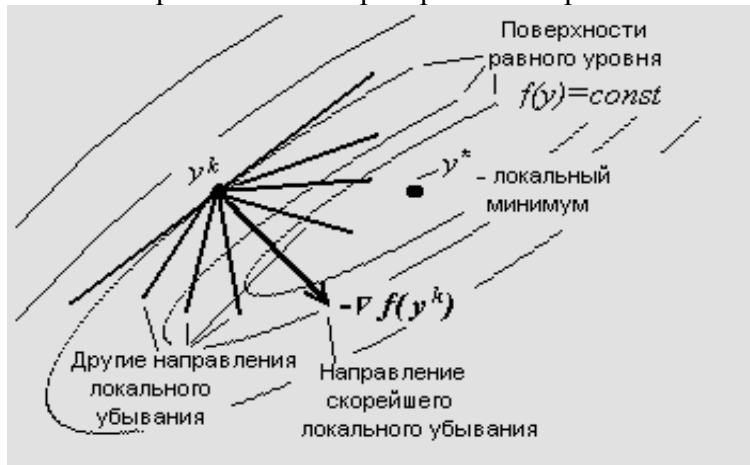


Рис. 4. 2. Отличие направления антиградиента от направления на локальный минимум

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 (4.2), то есть допустимая область $Q=D$ (общие принципы учета ограничений–неравенств рассмотрены в разделе 4.6). Применительно к этой задаче в методе наискорейшего градиентного поиска величина одномерного смещения x^k определяется из условия

$$f(y^k + x^k d^k) = \min \{f(y^k + x d^k) : x \geq 0, y^k + x d^k \in Q=D\}, \quad (4.15)$$

Таким образом, x — величина, определяющая смещение вдоль d^k , выбирается из условия достижения минимума функции f в области Q на луче $y^k + x d^k$, где $x \geq 0$. С вычислительной точки зрения правило (4.15) реализуется методом аккуратного одномерного поиска, описанным в разделе 4.2.4.

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 только локальная линейная модель. Выбор величины смещения по правилу (4.15) неявно предполагает нелинейность функции, особенно в задаче без ограничений.

В случае двух переменных правила поиска в методе наискорейшего градиентного поиска иллюстрирует рис. 4.3. Следует заметить, что направления поиска на двух последовательных шагах d^k и d^{k+1} взаимно ортогональны, если решение вспомогательной задачи (4.15) 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 $Q = R^N$) имеет следующую формулировку.

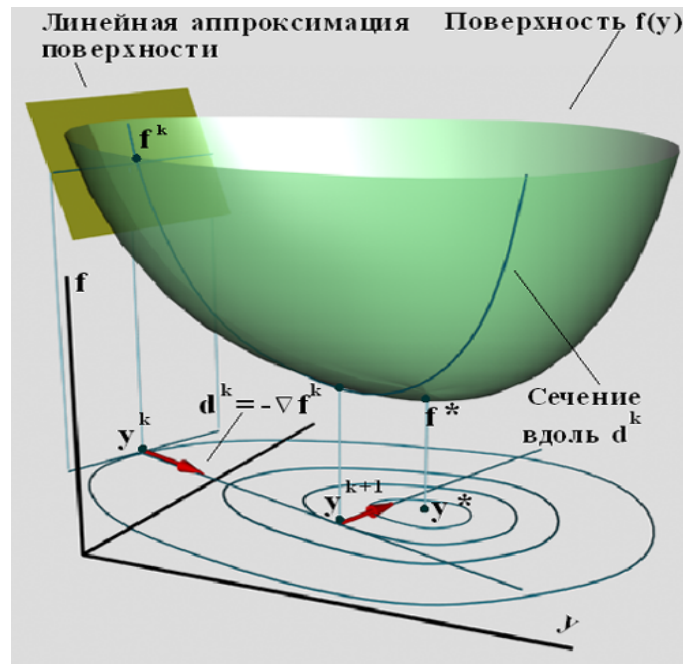


Рис. 4.3. Наискорейший градиентный поиск

Теорема 4.1. Пусть в задаче (4.4) функция $f(y)$ непрерывно дифференцируема, ограничена снизу и ее градиент удовлетворяет условию Липшица с некоторой константой L . 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 можно найти, например, в [1]. При определенных дополнительных предположения 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

$$f(y) = (y^T \Gamma y) / 2 + c^T y, \quad \Gamma^T = \Gamma \quad (4.16)$$

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 4.2. Для квадратичной функции (4.16) с симметричной положительно определенной матрицей метод наискорейшего градиентного поиска сходится со скоростью геометрической прогрессии со знаменателем, не превосходящим значения q из (4.16). 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 этого свойства можно найти в книге [8].

Из оценок теоремы 4.2 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.

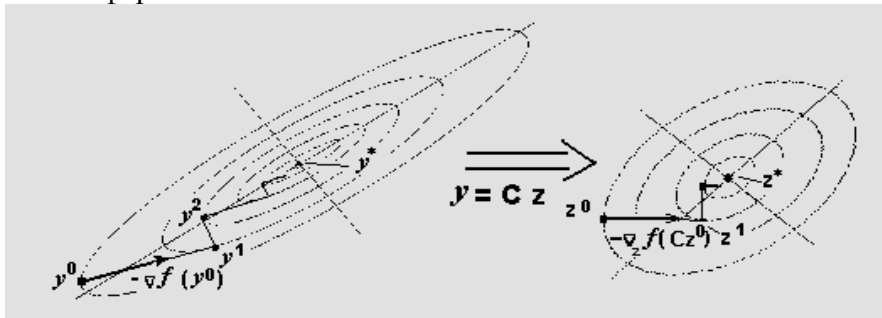


Рис. 4.4. Влияние масштабирования на скорость сходимости наискорейшего градиентного поиска

Если же поверхности равного уровня сильно вытянуты (рис.4.4), что соответствует $\lambda_{\min} \ll \lambda_{\max}$, то в (4.16) q будет близко к единице, и скорость сходимости к решению окажется 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 $z: y=Cz$ (C — матрица преобразования), то градиент функции в новых переменных z может быть вычислен как

$$\nabla_z f(Cz) = C^T \nabla f(y). \tag{4.17}$$

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

$$\bar{d}^k = C(-\nabla_z f(Cz^k)) = CC^T(-\nabla f(y^k)). \tag{4.18}$$

Возникающая матрица CC^T для коррекции направления антиградиента легко вычисляется для строго выпуклых квадратичных функций вида (4.16). 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 (4.16) совпадает со своей квадратичной аппроксимацией $P^k(y)$, построенной по измерениям $f^k, \nabla f^k, \Gamma_k^f$, выполненным в точке y^k

$$P^k(y) = (y - y^k)^T \Gamma_k^f (y - y^k) / 2 + (\nabla f^k, (y - y^k)) + f^k. \tag{4.19}$$

Условие, определяющее y^* — точку минимума для $P^k(y)$, примет вид

$$\nabla P^k(y^*) = \Gamma_k^f(y^* - y^k) + \nabla f^k = 0, \quad (4.20)$$

откуда $y^* - y^k = (\Gamma_k^f)^{-1}(-\nabla f^k)$. Сравнивая полученное направление с направлением (4.18), видим, что в качестве матрицы преобразования в (4.18) можно использовать $CC^T = (\Gamma_k^f)^{-1}$.

Если бы $f(y)$ была произвольной дважды непрерывно дифференцируемой функцией, то в (4.19) квадратичная аппроксимация $P^k(y)$ уже не совпадала бы с исходной функцией, а условие (4.20) 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 итерации

$$y^{k+1} = y^k + (\Gamma^f(y^k))^{-1}(-\nabla f(y^k)) \quad (4.21)$$

(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, выполняет 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 (4.21) 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.4.5.

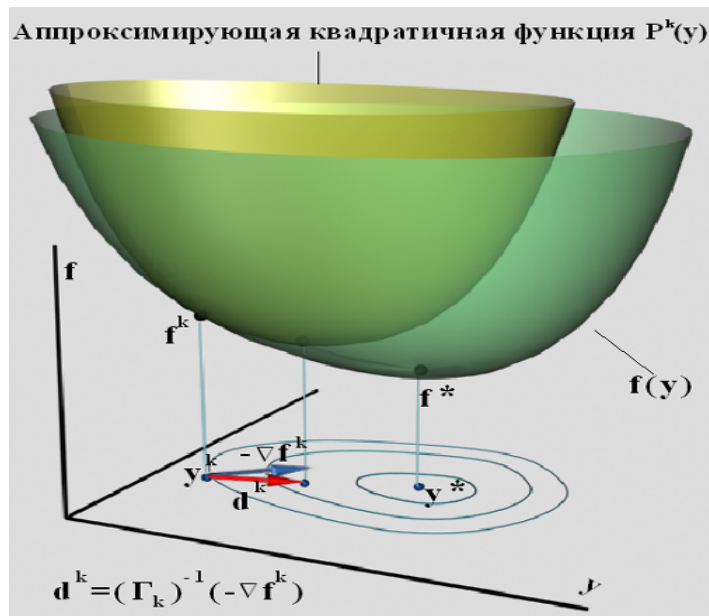


Рис. 4. 5. Выбор направления в методе Ньютона

Что можно в достаточно общем случае сказать о сходимости метода Ньютона?

Теорема 4.3. Для дважды непрерывно дифференцируемых функций 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 функция трижды непрерывно дифференцируема, то метод сходится квадратично.

Доказательство приведено, например, в [7]. Поясним терминологию.

Определение. *Линейной сходимостью* называют сходимость по закону геометрической прогрессии. Линейная сходимость характерна для метода наискорейшего градиентного поиска (теорема 4.2).

Определение. 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 $\|y^{k+m} - y^*\| \leq \alpha_{k+1} \dots \alpha_{k+m} \|y^k - y^*\|$.

Определение. Говорят, что метод *сходится квадратично*, если $\exists T > 0$, что $\|y^{k+1} - y^*\| \leq T \|y^k - y^*\|^2$ при $\|y^0 - y^*\| < \varepsilon$.

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 4.3 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 (4.21) 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 (4.21) неприменима. Если же $\Gamma^f(y^*)$ не вырождена, но не знакоположительна то, как следует

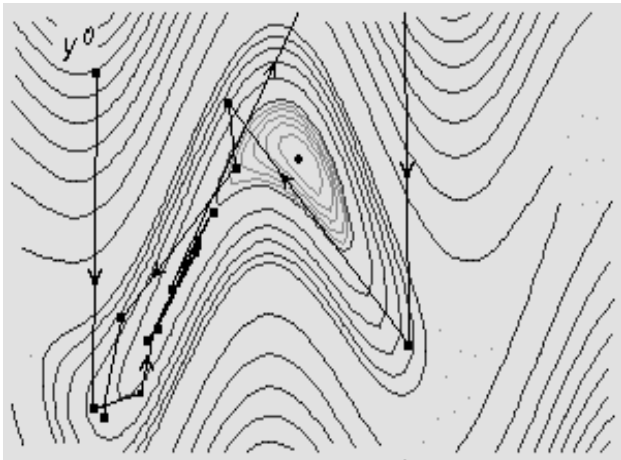


Рис. 4.6. Неожиданное поведение метода Ньютона

из приведенной выше теоремы, метод Ньютона может сходиться к стационарной точке функции f , не являющейся точкой минимума, а представляющей собой точку максимума или седловую точку.

Пример ситуации, в которой наблюдается сложное и несколько неожиданное поведение метода Ньютона, вызванное изменением знакоопределенности матриц вторых производных вдоль траектории поиска, показан на рис.4.6.

Все сказанное выше позволяет прийти к следующему выводу: для эффективного решения прикладных задач, характеризующихся плохим масштабированием, возможным вырождением и знаконеопределенностью матриц вторых производных, классические методы поиска локального экстремума являются малоприспособными. Необходимы методы, сочетающие сходимость из любой начальной точки с высокой скоростью сходимости вблизи решения и сохраняющие свои свойства в ситуациях, характерных для прикладных задач.

4.4. Методы локальной оптимизации, основанные на квадратичной модели поведения функций

Существует большая группа методов, в которых при выборе направления очередного шага используется предположение о том, что минимизируемая функция

хорошо приближается ее квадратичной аппроксимацией. Будем говорить, что такие методы *основаны на квадратичной модели* поведения функции.

Использование методом квадратичной модели вовсе не предполагает, что в нем обязательно явно используется матрица вторых производных функции или ее оценка. Есть методы, которые действительно явно используют матрицу Гессе. Есть методы, которые строят ее оценку, измеряя вдоль траектории поиска только вектор градиента, и, наконец, есть методы, которые явно матрицу Гессе не используют и не оценивают, хотя близость функции к квадратичной предполагают. В разделах 4.4.1, 4.4.2 рассмотрены методы каждого из этих видов.

4.4.1. 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.

4.4.1.1. Недостатки классического метода Ньютона. Анализ влияния регулировки величин одномерных смещений на свойства метода

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. 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 $\forall y \in R^N, z \in R^N$.

$$m\|z\|^2 \leq z^T \Gamma^f(y) z \leq M\|z\|^2 \quad (4.22)$$

Такие функции будут сильно выпуклы [1].

Можно показать, что условие (4.22) равносильно тому, что все собственные числа $\lambda_1(y), \dots, \lambda_N(y)$ матриц $\Gamma^f(y)$ лежат между m и M .



Замечание. Условие (4.22) 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 1. Любая функция $f(y)$ из класса $\Phi_{m,M}$ имеет единственный минимум y^* .

Свойство 2. Для функций $f(y)$ из класса $\Phi_{m,M}$ существует взаимосвязь между ошибкой по координате и ошибкой по значению функции, выражаемая соотношением

$$0,5m\|y-y^*\|^2 \leq f(y) - f^* \leq 0,5M\|y-y^*\|^2.$$

Свойство 3. Для функций $f(y)$ из класса $\Phi_{m,M}$ существует взаимосвязь между ошибкой по значению функции и нормой градиента

$$M(1+m/M)(f(y)-f^*) \leq \|\nabla f(y)\|^2.$$

Доказательство этих свойств можно найти в [1].

Свойства метода Ньютона с регулировкой шага на классе $\Phi_{m,M}$ определяются следующей теоремой [1]. Ее доказательство опирается на приведенные выше свойства.

Теорема 4.4 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.



Замечание. Если дополнительно потребовать от функции f существование непрерывных третьих производных, можно доказать, что метод будет сходиться квадратично.

Таким образом, модификация Ньютона–Рафсона расширяет область сходимости метода Ньютона.

4.4.1.2. 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. Таким образом, матрицу Γ_k использовать нельзя.

Основная идея, на основе которой выполняется модификация матриц, состоит в том, чтобы заменить матрицу 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

$$\begin{aligned} y^{k+1} &= y^k + x^k \bar{d}^k \\ \bar{d}^k &= (\bar{\Gamma}_k)^{-1} (-\nabla f(y^k)) \\ x^k &\in \Pi_1(\eta) \cap \Pi_2(\mu). \end{aligned} \quad (4.23)$$

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

$$\Gamma_k = \lambda_1 u^1 (u^1)^T + \dots + \lambda_N u^N (u^N)^T = U L U^T,$$

где матрица U составлена из вектор–столбцов u^1, \dots, u^N , 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

$$\bar{\lambda}_i = \begin{cases} \lambda_i, & \text{при } \lambda_i > \delta \\ \delta, & \text{при } \lambda_i \leq \delta, \end{cases} \quad (4.24)$$

где δ – малое положительное число. После этого полагается

$$\bar{\Gamma}_k = U \bar{L} U^T, \quad (4.25)$$

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 становится 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 (4.24)–(4.25), то произойдут качественные изменения в структуре аппроксимации. На рис.4.7 на примере пространства двух переменных показаны изменения в линиях равного уровня квадратичной аппроксимации $P^k(y)$ после замены знаконеопределенной матрицы Γ_k на положительно определенную $\bar{\Gamma}_k$. На этом рисунке видно, как направление метода Ньютона d^k , приводящее в стационарную точку поверхности $P^k(x)$ заменяется новым направлением \bar{d}^k .

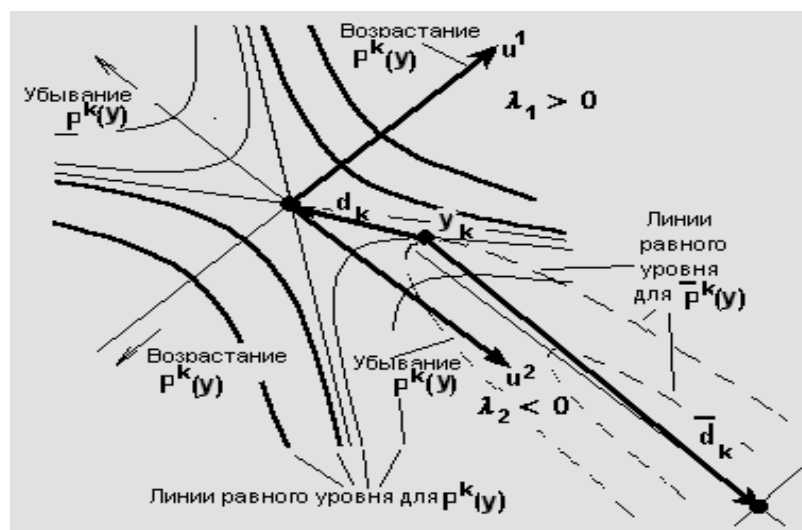


Рис. 4. 7. Изменение изолиний при замене матрицы

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. 4.8– 4.9.

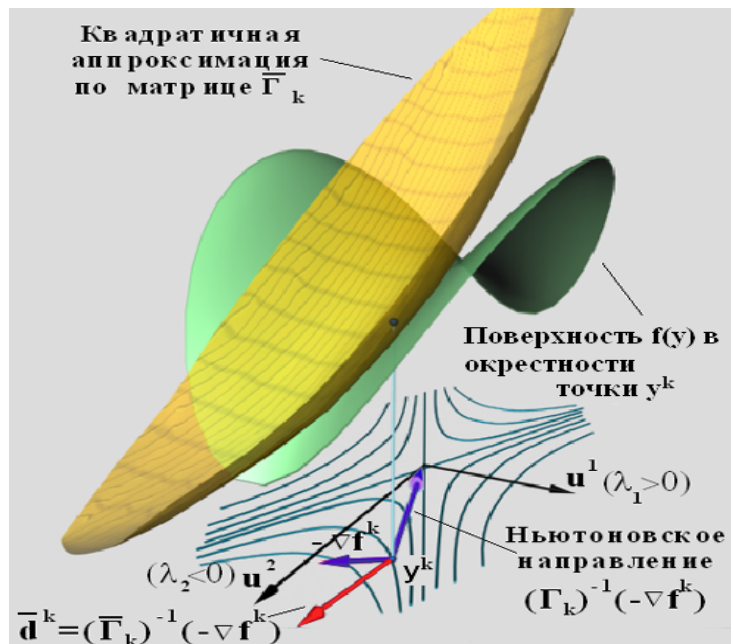


Рис. 4. 8. Изменение вида аппроксимирующей поверхности в случае знаконеопределенной матрицы Гессе

Поведение функции $f(x)$, представленное на рис.4.8, соответствует изолиниям, показанным на рис.4.7. 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 тоже время, в направлении u^2 отрицательная кривизна заменяется малой положительной кривизной.

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.4.9). Это приводит к замене направления Ньютона d^k на новое — \bar{d}^k , ориентированное на точку минимума измененной аппроксимации.

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 Холесского [2], вычислительная реализация которого проще и требует меньшего числа операций.

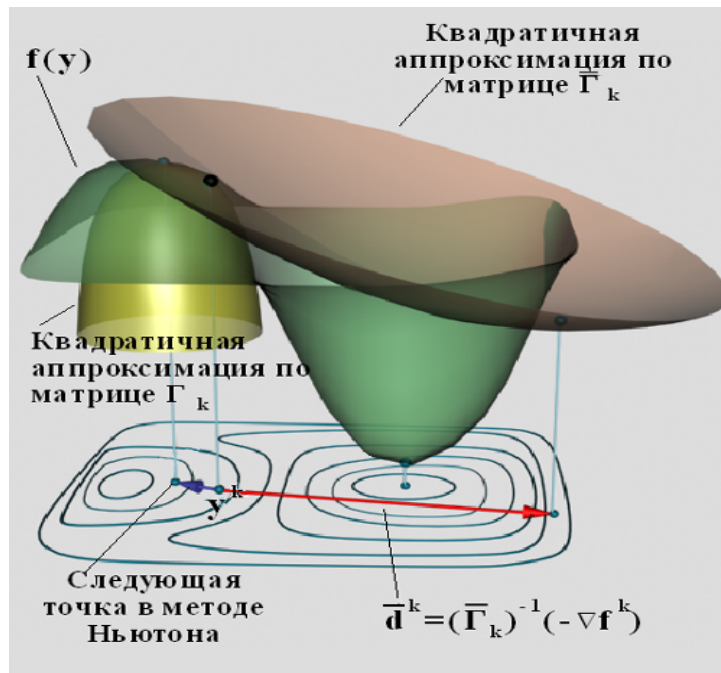


Рис. 4. 9. Влияние модификации матрицы при ее отрицательной определенности

В теории матриц известно, что для любой симметричной положительно определенной матрицы \mathbf{G} существует нижняя треугольная матрица \mathbf{L} с единичной диагональю и диагональная матрица \mathbf{D} с положительной диагональю, что справедливо разложение Холецкого $\mathbf{G} = \mathbf{LDL}^T$, т.е.

$$\begin{pmatrix} \gamma_{11} & \gamma_{12} & \dots & \gamma_{1N} \\ \gamma_{12} & \gamma_{22} & \dots & \gamma_{2N} \\ \cdot & \cdot & \dots & \cdot \\ \gamma_{1N} & \gamma_{2N} & \dots & \gamma_{NN} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{12} & 1 & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ l_{1N} & l_{2N} & \dots & 1 \end{pmatrix} \begin{pmatrix} d_{11} & 0 & \dots & 0 \\ 0 & d_{22} & \dots & 0 \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & d_{NN} \end{pmatrix} \begin{pmatrix} 1 & l_{12} & \dots & l_{1N} \\ 0 & 1 & \dots & l_{2N} \\ \cdot & \cdot & \dots & \cdot \\ 0 & 0 & \dots & 1 \end{pmatrix}, \quad (4.26)$$

где $d_{11}, \dots, d_{NN} > 0$.

Отсюда вытекает, что

$$\gamma_{11} = d_{11}, \gamma_{1j} = d_{11}l_{1j}, (j > 1) \quad (4.27)$$

$$\gamma_{ii} = d_{ii} + (l_{i1})^2 d_{11} + (l_{i2})^2 d_{22} + \dots + (l_{(i-1)(i-1)})^2 d_{(i-1)(i-1)}, (i = 2, \dots, N) \quad (4.28)$$

$$\gamma_{ij} = (l_{i1})(l_{1j})d_{11} + (l_{i2})(l_{2j})d_{22} + \dots + (l_{(i-1)(i-1)})(l_{(i-1)j})d_{(i-1)(i-1)} + (l_{ij})d_{ii}, \quad (4.29)$$

для $i < j \leq N$.

Из (4.27)–(4.29) легко получить формулы Холецкого для вычисления коэффициентов разложения Холецкого. Расчет их значений производится строка за строкой для матриц \mathbf{D} и \mathbf{L}^T . Порядок вычисления этих коэффициентов удобно пояснить следующей диаграммой:

$$d_{11} \Rightarrow l_{12}, l_{13}, l_{14}, \dots; d_{22} \Rightarrow l_{23}, l_{24}, \dots; d_{33} \Rightarrow l_{34}, l_{35}, \dots$$

Получим

$$d_{11} = \gamma_{11}, l_{1j} = (1/d_{11})\gamma_{1j}, (j = 2, \dots, N) \quad (4.30)$$

$$d_{ii} = \gamma_{ii} - \sum_{s=1}^{i-1} (l_{si})^2 d_{ss}, (i > 1)$$

$$l_{ij} = (\gamma_{ij} - \sum_{s=1}^{i-1} (l_{si})(l_{sj})d_{ss})/d_{ii}, (i > 1, j = i+1, \dots, N) \quad (4.31)$$

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

$$\bar{d}_{ij} \geq \delta > 0 \quad (4.32)$$

Это обеспечит положительность с «запасом» элементов модифицированной диагональной матрицы \bar{D} . Отметим, что условие (4.32) не может являться единственным условием модификации. Действительно, близость к нулю некоторых модифицированных элементов \bar{d}_{ij} при их дальнейшем использовании в (4.31) может привести к лавинообразному росту элементов l_{ij} при вычислениях. При обычном разложении Холецкого это невозможно, т.к. из (4.28) 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

$$d_{ss} \leq \beta^2; \quad ((\bar{l}_{si})^2 \bar{d}_{ss}) \leq \beta^2 \quad (s = 1, \dots, N, i > s) \quad (4.33)$$

Выполним необходимую модификацию за счет изменения только диагональных элементов матрицы Γ . Обозначим через Δ_i добавки к элементам γ_{ii} . Тогда согласно (4.32)–(4.33) должно выполняться:

$$\bar{d}_{ii} = \gamma_{ii} + \Delta_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss} \geq \delta > 0,$$

$$\max_{j=i+1, \dots, N} (\bar{l}_{ij})^2 \bar{d}_{ii} = \left(\left(\max_{j=i+1, \dots, N} \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right) \right)^2 / \left(\gamma_{ii} + \Delta_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss} \right) \right) \leq \beta^2.$$

Если обозначить

$$c_i^2 = \max_{j=i+1, \dots, N} \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right)^2 \quad (4.34)$$

$$\tilde{d}_{ii} = \gamma_{ii} - \sum_{s=1}^{i-1} (\bar{l}_{si})^2 \bar{d}_{ss}, \quad (4.35)$$

то

$$\Delta_i = \max \{0; c_i^2 / \beta^2 - \tilde{d}_{ii}; \delta - \tilde{d}_{ii}\}.$$

Это соответствует выбору

$$\bar{d}_{ii} = \max \{ \tilde{d}_{ii}; c_i^2 / \beta^2; \delta \} \quad (4.36)$$

$$\bar{l}_{ij} = \left(\gamma_{ij} - \sum_{s=1}^{i-1} (\bar{l}_{si})(\bar{l}_{sj}) \bar{d}_{ss} \right) / \bar{d}_{ii} \quad (4.37)$$

Элементы (4.36), (4.37) определяют модифицированные матрицы \bar{L}^T , \bar{D} . В качестве положительно определенной матрицы – приближения для Γ используется $\bar{\Gamma} = \bar{L} \bar{D} \bar{L}^T$.

Заметим, что сумма квадратов поправок к элементам матрицы Γ равна $\Delta_1^2 + \dots + \Delta_N^2$. Для уменьшения этой величины в [2] рекомендуется выбирать

$$\beta^2 = \max \{ \gamma^*; \xi / (N^2 - 1)^{1/2}; \varepsilon_M \}, \quad (4.38)$$

где ε_M — наименьшее положительное вещественное число в машинной арифметике, а

$$\gamma^* = \max \{ \gamma_{ii}; i = 1, \dots, N \}, \quad \xi = \max \{ |\gamma_{ij}|; 1 \leq i \leq N, 1 \leq j \leq N, i \neq j \}. \quad (4.39)$$

Приведем пошаговое описание метода Ньютона с регуляризацией шага и модификацией матрицы вторых производных на положительную определенность с использованием модифицированного преобразования Холесского.

ШАГ 0. Задаются начальная точка y^0 , параметры выбора коэффициента одномерного шага $0 < \mu < \eta < 1$, $0 < \sigma < 1$; 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 (4.38), (4.39). Строится модифицированное разложение 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 $x^k \in \Pi$ из (4.13). Определяется $y^{k+1} = y^k + x^k \bar{d}^k$

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.



Известно, что выполнение модифицированного разложения Холесского требует около $(1/6)N^3$ операций, а последующее определение \bar{d}^k требует числа операций порядка N^2 .

Замечание. 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 обладает свойством (4.22). Это следует из того, что в такой области модификация матрицы производных не будет, т.е. $\bar{\Gamma}_k = \Gamma_k$, а, следовательно, метод будет точно совпадать с методом Ньютона–Рафсона.

4.4.2. Методы первого порядка для гладких задач

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 Гессе и строятся на основе эвристических принципов. Методы третьей группы явно никаких матриц не строят, хотя неявно метрику пространства изменяют (методы сопряженных направлений).

4.4.2.1. 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 Гессеана 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 Гессеана 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 матрицу G_k , являющуюся оценкой кривизны поверхности $f(y)$ на траектории поиска. После выполнения $k-1$ испытания, G_k рассматривается как оценка матрицы вторых производных $\Gamma^f(y^k)$. Точка очередного измерения будет выбираться по правилу

$$y^{k+1} = y^k + x^k d^k, \quad (4.40)$$

$$d^k = (G_k)^{-1} (-\nabla f^k), \quad (4.41)$$

$$x^k \in \Pi.$$

Методы вида (4.40), (4.41) 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 матрицу G_k . Поэтому второе название этих методов – *методы переменной метрики*. При определенных условиях через N шагов матрица G_N будет близка к матрице $\Gamma^f(y^N)$.

Рассмотрим основные идеи, лежащие в основе построения оценки G_k . Рассмотрим случай, когда $f(y) = (y^T \Gamma y)/2 + (c, y) + b$ – квадратичная функция с симметричной положительно определенной матрицей Γ .

Пусть

$$\Delta^k = y^{k+1} - y^k,$$

$$z^k = \nabla f^{k+1} - \nabla f^k, \quad (4.42)$$

где $\nabla f^k = \nabla f(y^k)$, $\nabla f^{k+1} = \nabla f(y^{k+1})$. Тогда, очевидно, будет выполняться равенство

$$z^k = \nabla f(y^k + \Delta^k) - \nabla f(y^k) = \Gamma \Delta^k.$$

Потребуем, чтобы такому же условию удовлетворяла оценка G_{k+1} матрицы Γ , построенная по $(k+2)$ -м измерениям градиента. А именно, пусть выполняется требование

$$z^k = G_{k+1} \Delta^k. \quad (4.43)$$

Условие (4.43) называется *квазиньютоновским условием*.

Наложим дополнительные требования на матрицы оценок. Поскольку сама матрица \mathbf{G} симметрична, потребуем выполнения свойства симметрии от матрицы \mathbf{G}_{k+1} , положив

$$\mathbf{G}_{k+1} = (\mathbf{G}_{k+1})^T. \quad (4.44)$$

Будем определять ее новое значение путем коррекции предыдущей матрицы

$$\mathbf{G}_{k+1} = \mathbf{G}_k + \mathbf{U}_k.$$

где поправки \mathbf{U}^k строятся в виде матриц ранга 1 и находятся из условий (4.43), (4.44).

Эти условия определяют поправку неединственным образом. Простейший способ ее определения, предложенный Бройденом, состоит в том, чтобы составить ее из вектор–столбцов 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

$$\mathbf{G}_{k+1} = \mathbf{G}_k + (z^k - \mathbf{G}_k \Delta^k)(v^k)^T / ((v^k)^T \Delta^k). \quad (4.45)$$

Если положить в ней $v^k = (z^k - \mathbf{G}_k \Delta^k)$, то получим формулу Бройдена (B–формулу)

$$\mathbf{G}_{k+1} = \mathbf{G}_k + (z^k - \mathbf{G}_k \Delta^k)(z^k - \mathbf{G}_k \Delta^k)^T / ((z^k - \mathbf{G}_k \Delta^k)^T \Delta^k) \quad (4.46)$$

Существуют и другие способы оценивания. Например, можно несимметричную поправку в формуле (4.45) заменить похожей симметричной. Непосредственной проверкой можно убедиться (выполните эту проверку!), что для любого вектора v^k такого, что $(v^k)^T \Delta^k \neq 0$, соотношение

$$\begin{aligned} \mathbf{G}_{k+1} = & \mathbf{G}_k + ((z^k - \mathbf{G}_k \Delta^k)(v^k)^T + v^k(z^k - \mathbf{G}_k \Delta^k)^T) / ((v^k)^T \Delta^k) - \\ & - (v^k)(v^k)^T (z^k - \mathbf{G}_k \Delta^k)^T \Delta^k / ((v^k)^T \Delta^k)^2 \end{aligned} \quad (4.47)$$

дает оценочную матрицу, удовлетворяющую (4.43), (4.44).

Положив в (4.45) $v^k = z^k$, получим формулу Девидона–Флетчера–Пауэлла (DFP–формулу), представимую в следующем виде

$$\mathbf{G}_{k+1} = \mathbf{G}_k - \mathbf{G}_k \Delta^k (\Delta^k)^T \mathbf{G}_k / ((\Delta^k)^T \mathbf{G}_k \Delta^k) + z^k (z^k)^T / ((z^k)^T \Delta^k) + (\Delta^k)^T \mathbf{G}_k \Delta^k w^k (w^k)^T, \quad (4.48)$$

где

$$w^k = z^k / ((z^k)^T \Delta^k) - \mathbf{G}_k \Delta^k / ((\Delta^k)^T \mathbf{G}_k \Delta^k). \quad (4.49)$$

Поскольку в (4.48), (4.49) $w^k (w^k)^T$ – симметричная матрица и, как можно показать, $(w^k)^T \Delta^k = 0$ (проверьте это в качестве упражнения), то, в силу (4.43), последнее слагаемое в (4.48) можно отбросить. Получаемая таким образом укороченная формула называется формулой Бройдена–Флетчера–Гольдфарба–Шанно (BFGH–формула).

В приведенных итерационных соотношениях начальное значение выбирается в виде $\mathbf{G}_0 = \mathbf{E}$ (единичная матрица).

В теории оптимизации известно удивительное свойство описанных выше матричных оценок. Сформулируем его в виде теоремы.

Теорема 4.5. Для квадратичных функций f с положительно определенными матрицами вторых производных матрица \mathbf{G}_N , полученная с использованием процедур (4.40)–(4.42), а также B, DFP или BFGH–формул, будет совпадать с матрицей вторых производных \mathbf{G}^f функции f , а матрицы \mathbf{G}_k ($k \leq N$) будут симметричны и положительно определены.

Доказательство этого факта можно найти в [11].

Следствие. Для квадратичной функции f с положительно определенной матрицей Гессе квазиньютоновский метод (4.40)–(4.42), использующий матричные оценки,

построенные по B, DFP или BFGH–формулам попадет в минимум функции f на $(N+1)$ -м шаге, т.е. точка $y^{N+1} = y^*$.

Действительно, если теорема верна, то после $(N+1)$ -го измерения градиента в основных точках траектории поиска метод построит оценку $G_N = F^f$, поэтому следующий шаг будет выполнен, фактически, по правилу метода Ньютона и приведет в точку минимума f .

Построенные алгоритмы могут быть применены для достаточно произвольных функций f , не являющихся квадратичными. В этом случае обычно $G_N \neq F^f(y_N)$, поскольку матрица вторых производных не постоянна. После каждого N -го шага необходим повторный запуск метода из получаемой точки y^{N+1} . Кроме того, процесс поиска может привести к тому, что матрицы G_k могут оказаться вырожденными или знаконеопределенными. При этом направление шага d^k в (4.41) перестанет быть направлением убывания функции и величина смещения x^k в (4.40) окажется равной нулю. Простейший способ коррекции в этом случае состоит в замене направления, построенного по правилу (4.41), на обычное антиградиентное направление. Стратегия поиска в квазиньютоновских методах проиллюстрирована на рис.4.10.

На этом рисунке показан (пунктиром) возможный вид линий равного уровня для квадратичных аппроксимаций функции $f(y)$, построенных по матричным оценкам G_k . Поскольку $G_0 = E$, 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 пространстве с новой метрикой, связанной с матрицами G_k . Они могут сильно отличаться от направлений градиента в исходном пространстве.

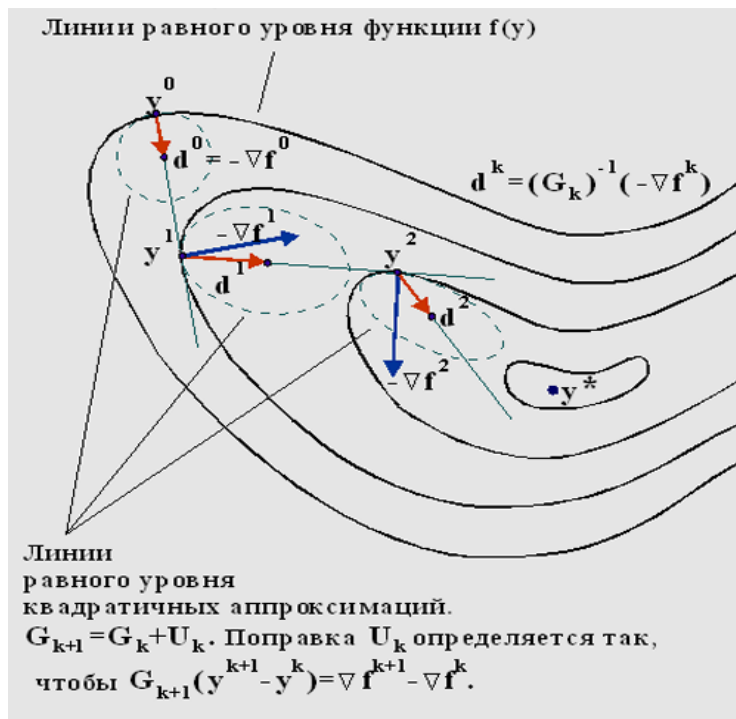


Рис. 4. 10. Направления поиска в квазиньютоновских методах

Замечание. При реализации алгоритма требуется проверка положительной определенности матриц G_k , а также вычисление направления поиска согласно

(4.41), 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 Холесского. Если при этом для диагональных элементов матрицы \mathbf{D} 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 $\mathbf{D}_k (\mathbf{L}_k)^T d^k = v$.

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 $\mathbf{G}_0 = \mathbf{E}$ и вычисляем $f^0 = f(y^0)$, $\nabla f^0 = \nabla f(y^0)$, $k = 0$.

ШАГ 2. Выполняем преобразование Холесского для матрицы \mathbf{G}_k . Если преобразование выполнить не удалось, полагаем $d^k = (-\nabla f^k)$ и переходим на шаг 4. В противном случае получаем $\mathbf{G}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^T$.

ШАГ 3. Определяем направление поиска $d^k = (\mathbf{G}_k)^{-1}(-\nabla f^k)$ путем решения двух систем с треугольными матрицами

$$\begin{aligned} \mathbf{L}_k v &= -\nabla f^k \\ \mathbf{D}_k (\mathbf{L}_k)^T d^k &= v \end{aligned} \quad (4.50)$$

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 6.

ШАГ 6. Производим вычисление матрицы \mathbf{G}_{k+1} по B -формуле (4.46), DFP -формуле (4.48) или по $BFGH$ -формуле. Переходим на шаг 2.

4.4.2.2. Модифицированные квазиньютоновские методы

В этих методах в случаях нарушения положительной определенности оценочных матриц \mathbf{G}_k вместо использования антиградиентного направления выполняется замена матрицы \mathbf{G}_k на близкую к ней положительно определенную матрицу $\bar{\mathbf{G}}_k$, построенную с использованием модифицированного преобразования Холесского (или другого подобного преобразования). Рис.4.11 показывает изменения вида квадратичной аппроксимации функции $f(y)$ в результате выполнения указанного преобразования для случаев различной знакоопределенности матрицы \mathbf{G}_k .

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 $\mathbf{G}_0 = \mathbf{E}$ и вычисляем $f^0 = f(y^0)$, $\nabla f^0 = \nabla f(y^0)$, $k = 0$.

ШАГ 2. Выполняем модифицированное преобразование Холесского для матрицы \mathbf{G}_k , получаем $\mathbf{G}_k \Rightarrow \bar{\mathbf{G}}_k = \bar{\mathbf{L}}_k \bar{\mathbf{D}}_k \bar{\mathbf{L}}_k^T$.

ШАГ 3. Определяем направление поиска $d^k = (\bar{\mathbf{G}}_k)^{-1}(-\nabla f^k)$ путем решения двух систем с треугольными матрицами

$$\bar{\mathbf{L}}_k v = -\nabla f^k$$

$$\bar{D}_k (\bar{L}_k)^T d^k = v$$

ШАГ 4. Определяем $x^k \in \Pi$ с помощью алгоритма одномерного шага, вычисляем

$$y^{k+1} = y^k + x^k d^k,$$

$$f^{k+1} = f(y^{k+1}), \nabla f^{k+1} = \nabla f(y^{k+1}),$$

$$\Delta_k = y^{k+1} - y^k, z^k = \nabla f^{k+1} - \nabla f^k.$$

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 (4.46), DFP -формуле (4.48) или по $BFGH$ -формуле. Переходим на шаг 2.

Практический опыт показывает, что для широкого класса гладких задач описанные алгоритмы достаточно экономичны по числу шагов.

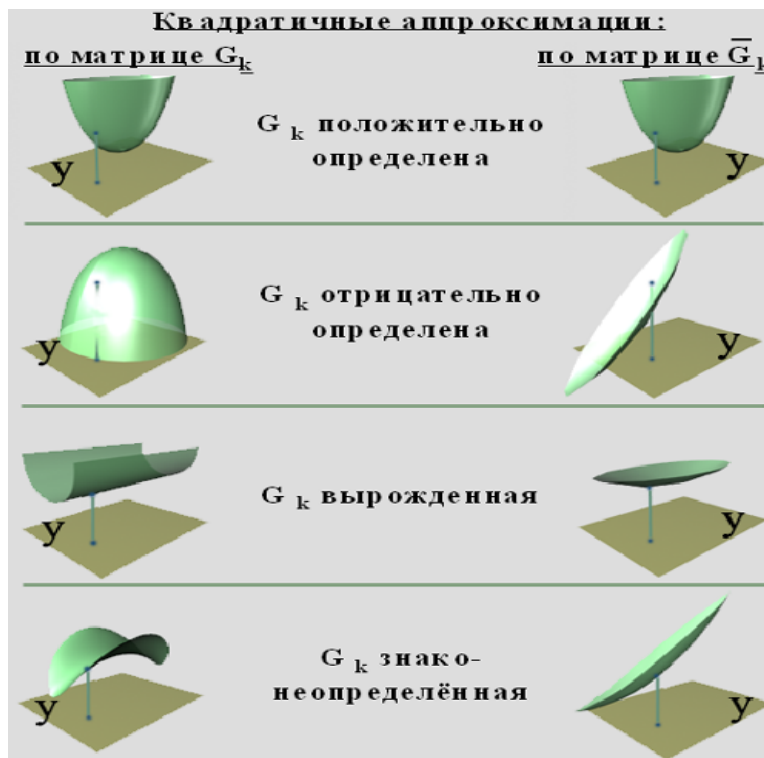


Рис. 4. 11. Влияние модификации матрицы G_k на вид аппроксимирующей поверхности

4.4.2.3. Эвристические методы коррекции метрики пространства поиска. R-алгоритмы Шора растяжения пространства

Автором этой группы методов является украинский математик Шор Н.З. Предложенные им методы основаны на эвристическом подборе матрицы преобразования пространства. Преобразования 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 ($e = 1, 2, \dots, k$), $\|\xi^e\| = 1$.

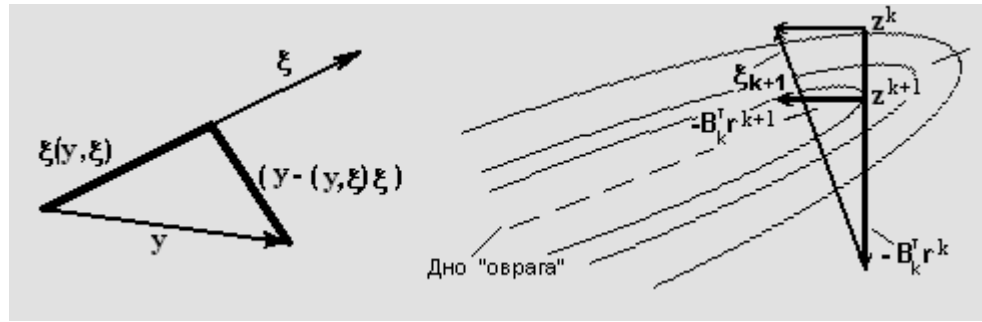


Рис. 4.12. Стратегия построения матрицы растяжения пространства

Нетрудно увидеть (рис. 4.12), что

$$R_\beta(\xi) y = (y - (y, \xi)\xi) + \beta(y, \xi)\xi = (E + (\beta - 1)\xi\xi^T)y. \quad (4.51)$$

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 . 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 (рис.4.13).

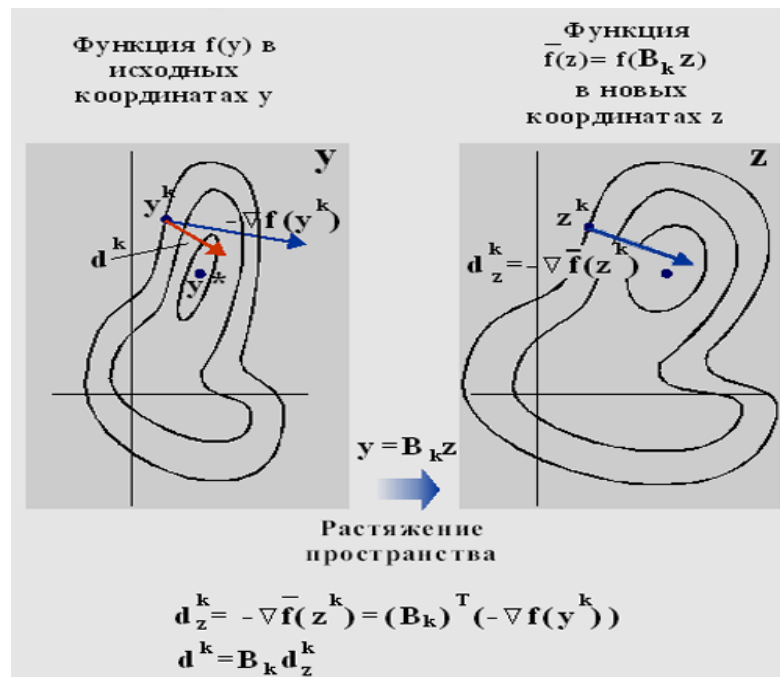


Рис. 4.13. Выбор направления в методе растяжения пространства

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.

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

$$\begin{aligned} y^{k+1} &= y^k + x^k B_k (B_k)^T (-r^k), \\ f^{k+1} &= f(y^{k+1}), \quad r^{k+1} = \nabla f(y^{k+1}). \end{aligned} \quad (4.52)$$

ШАГ 3. Если $\|r^{k+1}\| < \varepsilon$, то выполняется останов метода поиска, иначе переходим к шагу 4.

ШАГ 4. Выбирается направление дополнительного растяжения

$$\xi^{k+1} = (B_k)^T (r^{k+1} - r^k) \text{ и выполняется его нормировка } \xi^{k+1} := \xi^{k+1} / \|\xi^{k+1}\|.$$

ШАГ 5. Пересчитывается матрица преобразования с учетом растяжения пространства в α раз вдоль ξ^{k+1} :

$$B_{k+1} = B_k R_{1/\alpha}(\xi^{k+1}). \quad (4.53)$$

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.

4.4.2.4. Сопряженные направления и их свойства

Построение методов *сопряженных направлений* основано на квадратичной модели поведения минимизируемой функции. Предположим, что $f(y)$ — квадратичная функция (4.15) с положительно определенной матрицей.

Определение. Система линейно-независимых векторов p^0, p^1, \dots, p^{N-1} для симметричной матрицы Γ называется Γ -сопряженной, если

$$\forall i=1, \dots, N; j=1, \dots, N; i \neq j: (p^i, \Gamma p^j) = 0. \quad (4.54)$$

Определение. Пусть M — линейное многообразие, Γ — симметричная матрица, $x \neq 0$ и $x \in M$ и

$$\forall z \in M: (x, \Gamma z) = 0, \quad (4.55)$$

тогда вектор x называется Γ -сопряженным с многообразием M .

Можно легко доказать следующую лемму.

Лемма 4.1 Если p^0, p^1, \dots, p^{N-1} — все отличны от нуля, Γ — не только симметрична, но еще и положительно определена, тогда из (4.54) следует линейная независимость векторов p^0, p^1, \dots, p^{N-1} .

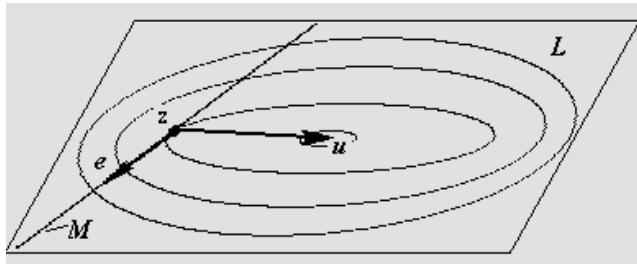


Замечание. В условиях леммы сопряженность означает ортогональность в смысле некоторого нового скалярного произведения.

Для построения методов, использующих сопряженные направления, чрезвычайно важным является свойство, определяемое следующей леммой.

Лемма 4.2 Пусть $f(y)$ — квадратичная функция вида (4.15) 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.4.14. 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 (4.55), следовательно $(u - z)$ будет Γ -сопряжен с многообразием M .

Рис. 4.14. Иллюстрация к лемме 4.2

Построим теперь вычислительные процедуры поиска минимума квадратичной функции $f(y)$, использующие Γ -сопряженные направления.

Определение. Поисковые процедуры вида (4.56), (4.57) называются методами сопряженных направлений.

$$y^{k+1} = y^k + x^k p^k \tag{4.56}$$

$$f(y^k + x^k p^k) = \min\{f(y^k + x p^k): -\infty < x < +\infty\}. \tag{4.57}$$

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 4.6. Пусть $f(y)$ — квадратичная функция вида (4.15) 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 (4.56), (4.57) 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 [8]. 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 4.2 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 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 можно придти 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.4.15).

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 $f(y)$.

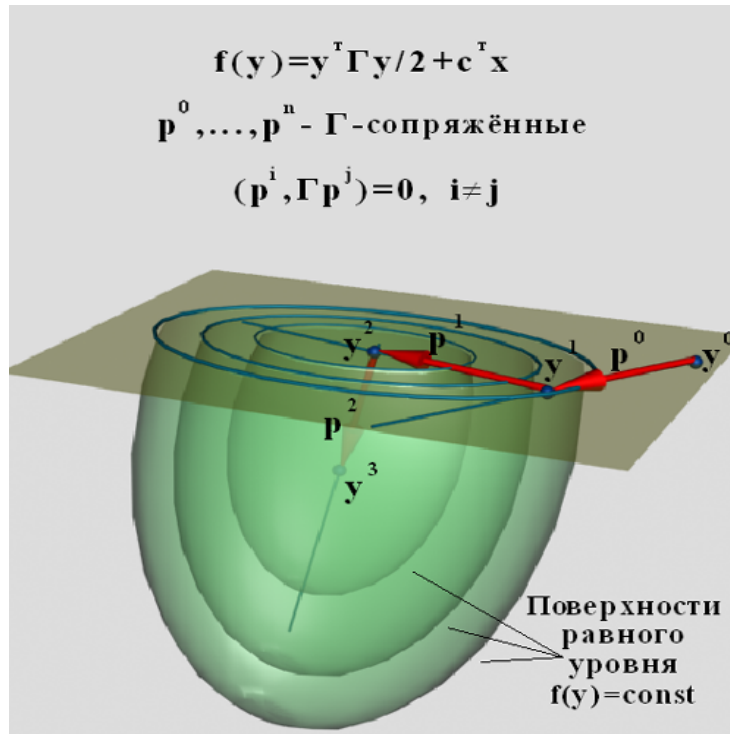


Рис. 4.15. Замечательное свойство сопряженных направлений

4.4.2.5. 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) [8]. Выберем

$$p^0 = -\nabla f^0, \quad \nabla f^0 = \nabla f(y^0). \tag{4.58}$$

Пусть векторы p^0, \dots, p^{k-1} построены. Положим

$$p^k = -\nabla f^k + \beta^{k-1} p^{k-1}, \quad \nabla f^k = \nabla f(y^k), \tag{4.59}$$

где y^k определяется условиями (4.56), (4.57). Подберем β^{k-1} из условия

$$(p^k, \Gamma p^{k-1}) = 0.$$

Получим

$$\beta^{k-1} = ((\nabla f^k)^T \Gamma p^{k-1}) / ((p^{k-1})^T \Gamma p^{k-1}) \tag{4.60}$$

Значение x^k , удовлетворяющее (4.57) для функции f , можно получить из условия $(\nabla f(y^k + x^k p^k), p^k) = 0$, 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

$$x^k = -((p^k)^T \nabla f^k) / ((p^k)^T \Gamma p^k) = -((\nabla f^k)^T \nabla f^k) / ((\nabla f^k)^T \Gamma p^k). \tag{4.61}$$

Для этого в числителе и знаменателе первой дроби необходимо выразить p^k из (4.59) 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 (4.56) на Γ , получим дополнительное соотношение

$$\nabla f^{k+1} = \nabla f^k + x^k \Gamma p^k. \tag{4.62}$$

Лемма 4.3. 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 (4.59), (4.62), (4.61), лемму можно доказать методом математической индукции [8].

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 (4.61) 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 (4.60) для вычисления коэффициента β^{k-1} вошла неизвестная матрица Γ . Это не является существенным, поскольку формула (4.60) может быть переписана в другом виде. Чтобы показать это, выразим Γp^{k-1} в числителе (4.60) из (4.62), а p^{k-1} в знаменателе (4.60) из (4.59). 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 (4.62) и пользуясь ортогональностью ∇f^k и ∇f^{k-1} , окончательно получим

$$\beta_{k-1} = \|\nabla f^k\|^2 / \|\nabla f^{k-1}\|^2 \quad (4.63)$$

Построенный метод определяет минимум любой квадратичной функции с положительно определенной матрицей Гессе за N шагов. Заметим, 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 $p^k = -\nabla f^k$ и полагаем $k = 0$. Иначе 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 (4.63) и p^k по формуле (4.59). Переходим на шаг 2.

Что известно о скорости сходимости построенного метода? Можно показать [1], что для функций из класса $\Phi_{m,M}$, описанного в (4.22), 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.

4.5. 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 не только может 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 [3], метод деформируемого многогранника Нелдера-Мида [12] и его модификация – комплексный метод Бокса [13]. 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.

4.5.1. 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 противоположного 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 $a > 0$

$$y^* = \bar{y} + a(\bar{y} - y^h)$$

и вычисляем $f^* = f(y^*)$.

ШАГ 6. 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;

если $f^* > f^e$, то переходим на шаг 7.

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 $y^h := \hat{y}$, $f^h := \hat{f}$ 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).

Преобразования симплекса в пространстве R^2 при операциях отражения, растяжения и сжатия показаны на рис.4.16.

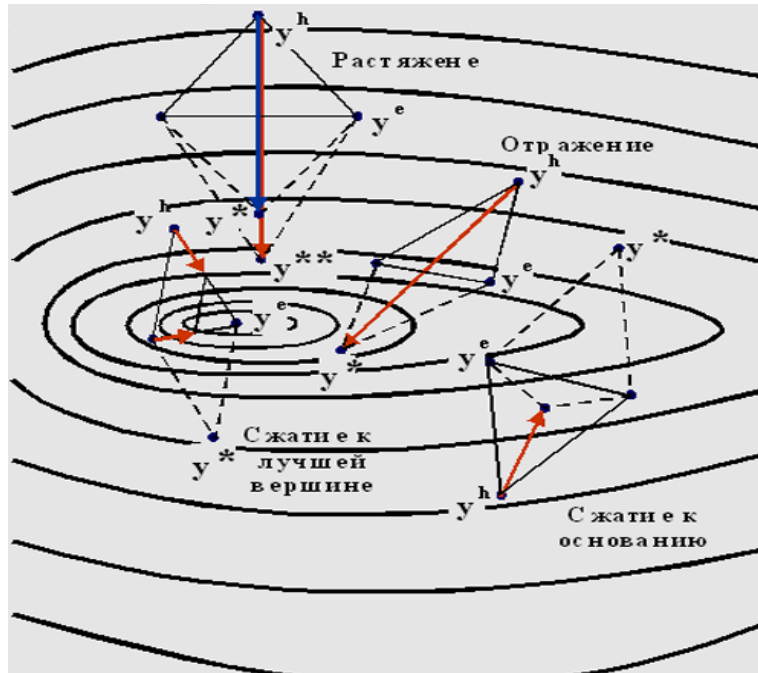


Рис. 4.16. Типовые операции с симплексом в методе Нелдера–Мида

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.

4.5.2. Метод Хука-Дживса

В этом разделе приводится краткое описание метода Хука-Дживса, который был специально разработан именно для задач с оврагами [3]. 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 $\bar{y} := \bar{y} - he^i$.

На рис.4.17 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.4.17.

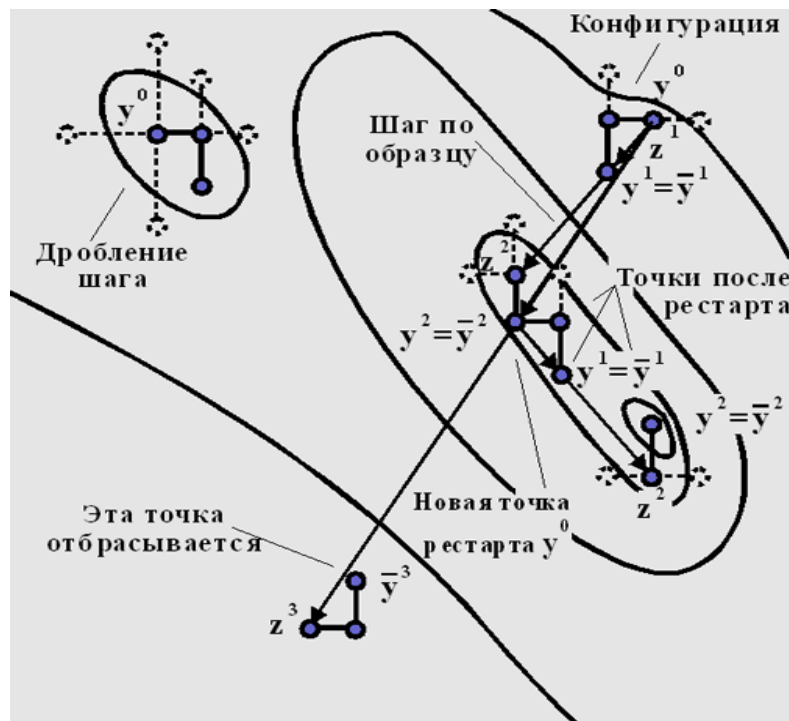


Рис. 4.17. Перемещения точки в методе Хука-Дживса

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.

4.6. Особенности применения методов локального поиска при двусторонних ограничениях на переменные

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 (4.3) в задачах могут присутствовать функциональные ограничения общего вида (4.2). В дальнейшем мы будем разделять два вида этих ограничений.

Общие ограничения, определяемые через функции ограничений, проще всего учесть с помощью одного из общих методов учета ограничений (такого, например, как метод внутренних [14] и внешних штрафов [6, 10], метод модифицированных функций Лагранжа [15] и другие [6, 16]), сводящего задачу с функциональными ограничениями к серии задач без функциональных ограничений. Такой подход позволяет применить к решению задач с ограничениями методы, ранее разработанные для задач без ограничений.

Двусторонние ограничения на переменные также можно было бы учесть на основе такой же методики. Однако их вид настолько прост, с одной стороны, а, с другой стороны, специфичен, что наиболее правильным решением является использование специальных методов для их учета. Это неизбежно приводит к пересмотру ранее рассмотренных алгоритмов и созданию их модификаций для задач с двусторонними ограничениями. Заметим, что двусторонние ограничения являются частным случаем линейных ограничений, специальные алгоритмы учета которых рассмотрены, например, в [2].

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.

4.6.1. Особенности учета двусторонних ограничений на переменные в методах гладкой оптимизации

Характерными чертами многих методов гладкой оптимизации (для задач без ограничений) являются:

- выполнение рестартов из последней достигнутой точки через определенное число шагов, которое зависит от размерности пространства поиска (например, методы квазиньютоновского типа, метод сопряженных градиентов);
- выполнение одномерного поиска в выбранном направлении или же перемещения в выбранном направлении с заранее заданным коэффициентом величины шага.

Наличие двусторонних ограничений будет оказывать влияние на реализацию каждого из этих двух процессов. А именно, процессы одномерных перемещений могут выводить на фрагменты границы области D из (4.3), определяемой ограничениями на переменные. После выхода на границу поиск должен продолжаться на линейном 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, выполняющие рестарты, должны будут производить их чаще, чем при поиске во всем пространстве. Это связано с тем, что периодичность рестартов определяется размерностью многообразия, на котором выполняется поиск. Еще одним важным моментом является правило возврата процесса поиска с текущего многообразия на многообразие (или в пространство) более высокой размерности в том случае, когда это приводит к уменьшению значения функции.

Дадим более точное описание правил выполнения всех 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.

$$J_a := J_a \cup \{i: 1 \leq i \leq N; y_i^k = a_i; d_i^{k-1} < 0\} \quad (4.64)$$

$$J_b := J_b \cup \{i: 1 \leq i \leq N; y_i^k = b_i; d_i^{k-1} > 0\} \quad (4.65)$$

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 незафиксированных компонент. 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 Z_k , используя их как вектор-столбцы

$$Z_k = (e^{j^1}, \dots, e^{j^n}).$$

Введем вектор переменных u для возникшего линейного подпространства R^n . Переход процесса поиска на линейное многообразие в исходном пространстве R^N

равносильна замене переменных $y = y^k + Z_k u$ с переходом к поиску в пространстве переменных u .

Вычисляя в старых переменных значения ∇f^k и Γ_k легко пересчитать их в соответствующие значения в новых переменных, используя известные соотношения

$$\nabla_u f^k = (Z_k)^T \nabla f^k; \quad \Gamma_{uk} = (Z_k)^T \Gamma_k Z_k$$

Нужно обратить внимание на то, что в пространстве новых переменных методы гладкой оптимизации должны быть запущены заново из точки $u^0=0$, соответствующей текущей точке поиска y^k . Заметим также, что если используемый метод включает рестарты, то при поиске минимума по переменным u 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 u происходит по обычным правилам, характерным для выбранного метода. Однако, при реализации этих правил необходимо все вектора и матрицы использовать для размерности n нового пространства, т.е., в частности, вместо значений ∇f^k и Γ_k необходимо использовать $\nabla_u f^k$ и Γ_{uk} .

Работа методов в пространстве переменных u 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 u выбрано направление поиска d^k_u . Тогда выполняется пересчет направления d^k_u в направление d^k в исходном пространстве переменных:

$$d^k = Z_k d^k_u.$$

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 разделе 4.2.4, 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, повторно применив правила их коррекции (4.64), (4.65) 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 x длины шага уменьшается таким образом, чтобы точка $y^{k+1} = y^k + x^k d^k$ 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 u не привело к выходу точки y^{k+1} на новые фрагменты границы. В этом случае необходимо проверить условия возврата к многообразию или пространству более высокой размерности за счет исключения из множеств J_a или J_b номеров части переменных. Правила исключения следующие:

$$J_a := J_a \setminus \{i \in J_a: \partial f(y^k)/\partial y_i < 0\} \quad (4.66)$$

$$J_b := J_b \setminus \{i \in J_b: \partial f(y^k)/\partial y_i > 0\}. \quad (4.67)$$

Условия исключения переменной в (4.66), (4.67) 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 (4.66), (4.67) 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.

4.6.2. 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 возможна специальная модификация этого метода для случая ограничений, рассмотренная, например, в [17]. Учет двусторонних ограничений на переменные в этом методе можно выполнить с использованием общего метода внешнего штрафа, который рассматривается в следующем разделе.

4.7. Учет ограничений общего вида на основе метода штрафов

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. Именно этот частный, но чрезвычайно важный случай был рассмотрен в предыдущем разделе.

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. В наиболее наглядной форме эта идея реализована в методе внешних штрафных функций, рассмотренных в следующем разделе.

4.7.1. Метод внешнего штрафа. Общие условия сходимости

Рассмотрим задачу в постановке (4.1)–(4.3). Через Q в ней обозначено множество допустимых точек, а через D — множество точек, удовлетворяющих двусторонним ограничениям на переменные ($Q \subseteq D$). Построим непрерывную в D функцию $H(y)$ так, чтобы она была положительной тогда и только тогда, когда $y \notin Q$, а в остальных точках обращалась в ноль:

$$H(y) > 0 \text{ при } y \notin Q \text{ и } H(y) = 0 \text{ при } y \in Q. \quad (4.68)$$

Функция такого вида называется *функцией штрафа*. Она определяет величину штрафа, накладываемого за выход точки y из допустимой области Q . В самой допустимой области штраф равен нулю.

При построении вспомогательных задач штрафная добавка добавляется к минимизируемой функции для того, чтобы препятствовать выходу точки поиска из допустимой области. В последующем штрафная добавка формируется в виде произведения функции штрафа на положительный числовой множитель, называемый *коэффициентом штрафа*. Использование коэффициента штрафа позволяет легко изменять (увеличивать) значение штрафной добавки, что, в общем случае, необходимо для существования подпоследовательностей решений возникающих вспомогательных задач, сходящихся к решениям исходной задачи.

Для реализации метода штрафов необходимо построить функцию штрафа $H(y)$ с указанными выше свойствами. Поскольку геометрическая структура области Q в общем случае неизвестна, принадлежность точек y к этой области определяется через значения функций ограничений в этих точках: если хотя бы одно ограничение–неравенство $g_i(y) \leq g_i^+$ нарушено, точка не принадлежит Q . Можно использовать функции штрафа разного вида [6,10]. Одна из распространенных форм — степенная функция штрафа, имеющая вид

$$H(y) = \sum_{i=1}^m (\max\{c_i (g_i(y) - g_i^+); 0\})^p. \quad (4.69)$$

Здесь $c_i > 0$ — нормирующие коэффициенты, выравнивающие диапазоны значений функций ограничений и приводящие их к одной размерности, p — параметр функции штрафа, влияющий на порядок ее гладкости.

Образует вспомогательную задачу со штрафом

$$S_\gamma(y) = f(y) + \gamma H(y), \quad (4.70)$$

$$S_\gamma(y) \rightarrow \min, y \in D, \quad (4.71)$$

$$D = \{y: a_i \leq y_i \leq b_i, i=1, \dots, N\}.$$

Здесь $\gamma > 0$ — коэффициент штрафа.

Почти очевидно, что решения конкретной задачи со штрафом в общем случае не будут являться решениями исходной задачи. Для их определения придется рассматривать не одну задачу вида (4.70), (4.71), а последовательность таких задач с возрастающими значениями коэффициента штрафа γ .

Такой подход к решению задач с ограничениями общего вида называют *методом внешнего штрафа*. Еще раз отметим, что он заключается в замене решения исходной задачи с функциональными ограничениями (4.1)–(4.3) на решение последовательности задач вида (4.70), (4.71) без функциональных ограничений, полученных для возрастающей последовательности значений *коэффициента штрафа* γ .

Замечание. Если функция штрафа (4.69) является гладкой, то нельзя гарантировать, что решение задачи со штрафом (4.71) при каком-либо конечном γ будут совпадать с решением исходной задачи (4.1)–(4.3). 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 Q в направлении d , выводящим из множества Q и не выводящим из D . 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 Q 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 .

Следует обратить внимание на то, что порядок гладкости функции степенного штрафа (4.69) легко регулируется за счет выбора соответствующего значения показателя степени p .

При $p=1$ функция штрафа (4.69) не является дифференцируемой по переменным g_i , при $p > 1$ она становится непрерывно дифференцируемой, а при $p > 2$ — дважды непрерывно дифференцируемой.

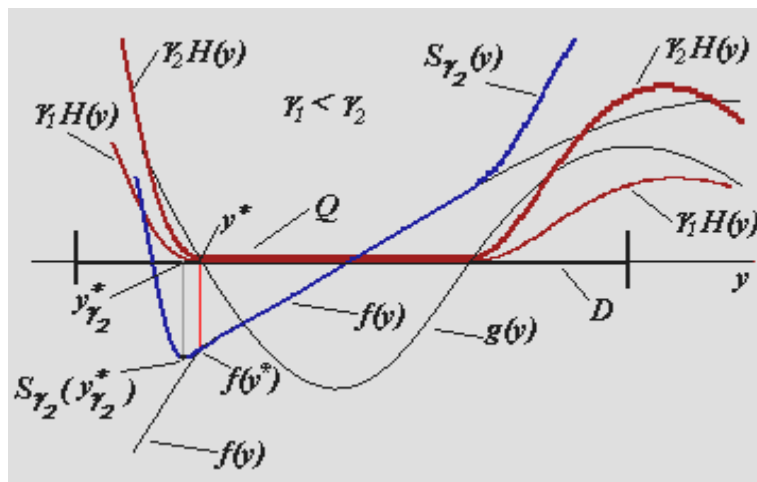


Рис. 4.18. Пример поведения функции штрафной задачи

На рис. 4.18 показано поведение функций задачи со штрафом в случае одного переменного и одного ограничения $g(y) \leq 0$ при показателе степени в штрафе $p=2$. Можно видеть изменение функции штрафа при увеличении коэффициента γ , а также вид функции штрафной задачи $S_\gamma(y)$ при одном из значений коэффициента штрафа.

Следует обратить внимание на то, что в задаче, представленной на рисунке, 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 $\gamma \rightarrow \infty$.

Теорема 4.7. Пусть функции $f(y)$ и $H(y)$ непрерывны в компактной области D и существуют глобальные минимумы $y^*_{\gamma_k}$ задач со штрафом (4.70), (4.71) при $\gamma = \gamma_k$, $k=1,2,\dots$. Тогда, если $\gamma_k \rightarrow \infty$ при $k \rightarrow \infty$, то все предельные точки последовательности $y^*_{\gamma_k}$ — решений задач со штрафом будут являться глобальными минимумами исходной задачи. При этом

$$\lim_{k \rightarrow \infty} f(y^*_{\gamma_k}(\varepsilon_k)) = f(y^*), \quad \lim_{k \rightarrow \infty} \rho(y^*_{\gamma_k}(\varepsilon_k), Y^*) = 0,$$

где Y^* — множество решений исходной задачи.

Доказательство. Пусть y^* — один из глобальных минимумов исходной задачи. Поскольку он принадлежит области D , по которой берется минимум в задаче со штрафом, то $S_{\gamma_k}(y^*) \geq S_{\gamma_k}(y^*_{\gamma_k}) = f(y^*_{\gamma_k}) + \gamma_k H(y^*_{\gamma_k})$.

Точки $y^*_{\gamma_k}$ принадлежат области D , которая в рассматриваемой постановке является компактом. Возьмем произвольную предельную точку y^*_∞ этой последовательности. Докажем, что эта точка допустима. Для этого разделим полученное выше неравенство на γ_k и перейдем к пределу на подпоследовательности $k=k_t$, для которой $y^*_{\gamma_{k_t}} \rightarrow y^*_\infty$ при $t \rightarrow \infty$. Поскольку коэффициент штрафа стремится к бесконечности, а функции f , S и H непрерывны на компакте D , а значит — ограничены, то в пределе получим, что $0 \geq H(y^*_\infty)$. В силу неотрицательности функции штрафа видим, что $H(y^*_\infty) = 0$, а следовательно $y^*_\infty \in D$, т. е. предельная точка допустима.

Докажем теперь, что эта предельная точка является глобальным минимумом исходной задачи. Используя еще раз записанное выше неравенство, усиливая его, получим: $f(y^*) \geq S_{\gamma_k}(y^*) \geq S_{\gamma_k}(y^*_{\gamma_k}) = f(y^*_{\gamma_k}) + \gamma_k H(y^*_{\gamma_k}) \geq f(y^*_{\gamma_k})$.

Переходя к пределу на той же подпоследовательности получим, что $f(y^*) \geq f(y^*_\infty)$. С другой стороны, поскольку из только что доказанного вытекает допустимость предельной точки y^*_∞ , обязательно выполняется обратное неравенство. Следовательно, y^*_∞ является глобальным минимумом исходной задачи.

Осталось обосновать два предельных соотношения, приведенных в конце теоремы. Следует обратить внимание на отличие этих новых утверждений от уже доказанных, состоящее в том, что них рассматриваются пределы последовательностей в целом, а не их сходящихся подпоследовательностей. Докажем первое из указанных предельных соотношений, имея ввиду, что второе соотношение доказывается с использованием аналогичных рассуждений.

Предположим, что соотношение $\lim_{k \rightarrow \infty} f(y^*_{\gamma_k}(\varepsilon_k)) = f(y^*)$ не верно. Тогда найдется подпоследовательность $k=k_t$, для которой при некотором $\varepsilon > 0$ начиная с $t > \bar{t}$ всегда будет выполняться $|f(y^*_{\gamma_{k_t}}) - f(y^*)| > \varepsilon$. Однако подпоследовательность $y^*_{\gamma_{k_t}}$, $k=k_t$ принадлежит компакту и, следовательно, имеет сходящуюся к y^* подпоследовательность, что противоречит последнему неравенству. Теорема доказана.

Доказанная теорема является теоретической основой метода штрафов, но не может служить в качестве его обоснования при вычислениях. Дело в том, что при численной оптимизации решения задач со штрафом никогда не являются точными. Это может привести к нарушению сходимости. Во всяком случае, необходимо дополнительное обоснование метода с учетом наличия ошибок в решениях штрафных задач.

Теорема 4.8. Пусть функции $f(y)$ и $H(y)$ непрерывны в компактной области D и существуют глобальные минимумы $y^*_{\gamma_k}$ задач со штрафом (4.70), (4.71) при $\gamma = \gamma_k$, $k=1,2,\dots$. Пусть также $y^*_{\gamma_k}(\varepsilon_k)$ — оценки этих глобальных минимумов с точностью ε_k по значению функции штрафной задачи, т.е. $S_{\gamma_k}(y^*_{\gamma_k}(\varepsilon_k)) \leq S_{\gamma_k}(y^*_{\gamma_k}) + \varepsilon_k$. Тогда при $\gamma_k \rightarrow \infty$, $\varepsilon_k \rightarrow 0$, $k \rightarrow \infty$ все предельные точки последовательности $y^*_{\gamma_k}(\varepsilon_k)$ — приближенных решений задач со штрафом будут являться глобальными минимумами исходной задачи. При этом

$$\lim_{k \rightarrow \infty} f(y^*_{\gamma_k}(\varepsilon_k)) = f(y^*), \quad \lim_{k \rightarrow \infty} \rho(y^*_{\gamma_k}(\varepsilon_k), Y^*) = 0,$$

где Y^* — множество решений исходной задачи.

Доказательство легко проводится по той же схеме, что и доказательство предыдущей теоремы. Его также можно найти, например, в [6, 10].

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

$$G(y) = \max \{ \max \{ c_j (g_j(y) - g_j^+); 0 \}; j=1, \dots, m \}, \quad (4.72)$$

ОПИСАНИЕ АЛГОРИТМА настройки параметров метода штрафов.

ШАГ 0. Задаются: $\varepsilon_0 > 0$ — начальная точность решения штрафных задач, $\varepsilon > 0$ — требуемая точность решения штрафных задач, $\delta > 0$ — требуемая точность по ограничениям, $0 < \alpha < 1$ — ожидаемый коэффициент убывания невязки по ограничениям на шаге, $\beta > 1$ — коэффициент увеличения штрафа, $\beta_1 > 0$ — 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 (4.70), (4.71) 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, иначе — сразу переходим на шаг 3.

ШАГ 3. Решаем задачу со штрафом (4.70), (4.71) 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 $G_{k+1} < \beta G_k$, 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 4.8, 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.

4.7.2. 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 (4.71), необходимо изучить характерные особенности, имеющиеся в структуре $S_{\gamma_k}(y)$ — функций штрафных задач. В начале данного раздела 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 Q , будут возникать функции $S_{\gamma_k}(y)$ с сильно овражной структурой.

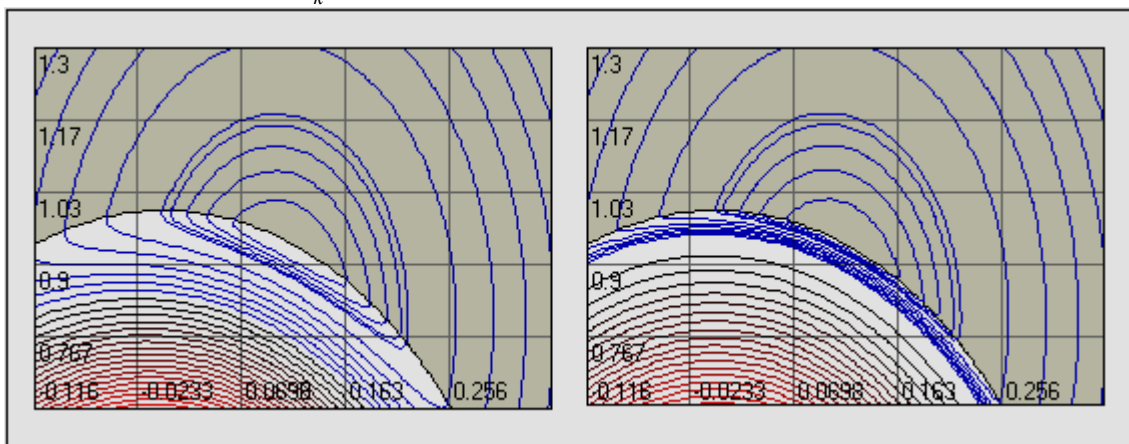


Рисунок 4.19 Увеличение овражности функции задачи со штрафом $S_{\gamma}(y)$ при возрастании коэффициента штрафа γ с 1 до 20

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, то, за счет больших значений коэффициента штрафа, функция $S_{\gamma_k}(y)$ будет быстро изменяться. Таким образом, у этой функции наблюдается «овраг».

На рис.4.19 приведен пример, показывающий изменение изолиний функции $S_{\gamma}(y)$ штрафной задачи вида (4.69)–(4.71), 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 $p=2$. Допустимая область выделена более темным цветом. На рис. 4.20 для этого же примера показана пространственная структура функции задачи со штрафом $S_{\gamma}(y)$ при небольшом значении коэффициента штрафа $\gamma=1$. По отношению к

изображению изолиний на рис. 4.19, пространственное изображение развернуто примерно на 80° .

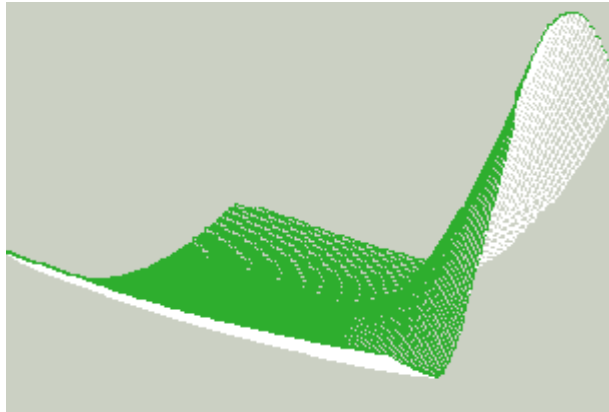


Рис. 4. 20. Пространственная структура функции $S_\gamma(y)$ при $\gamma=1$

Вторая особенность связана с тем, что функция штрафа $H(y)$ из (4.68), (4.69), 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 (4.69) при показателе степени $p \leq 1$ на границе нарушения ограничений все частные производные, в общем случае, будут разрывны, а при $1 < p \leq 2$ 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 Q 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 (4.1)–(4.3) лежит на границе области Q) 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 p в функции штрафа [6,10]. Эти оценки приведены ниже в теореме 4.9.

Третья особенность функции $S_{\gamma_k}(y)$ связана с тем, что порядок ее роста при отклонениях от границы области Q может быть существенно различен, в зависимости от того, происходит отклонение внутрь области Q 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).



Таким образом, функции штрафных задач обладают характерными особенностями, ухудшающими поисковые возможности методов локальной оптимизации.

Приведем формулировку упомянутой ранее теоремы об оценке скорости сходимости метода штрафов при использовании функции степенного штрафа (4.69).

Теорема 4.9 Пусть функция $f(y)$ липшицева в метрике $\rho(\cdot, \cdot)$ с константой L на компакте D , а функции ограничений–неравенств непрерывны на D и для них существуют такие $\delta > 0$ и $\alpha > 0$, что для любого y из δ -окрестности Q , т.е. для $y \in (O_\delta(Q) \setminus Q) \cap D$, гарантируется что функции ограничений будут возрастать не медленнее, чем линейные по $\rho = \rho(y, Q)$ функции: $\max\{c_j g_j(y) : j=1, \dots, m\} \geq \alpha \rho(y, Q)$.

Пусть также задачи со штрафом (4.70), (4.71) решаются точно, и используется функция степенного штрафа (4.69). Тогда справедливы следующие оценки для ошибки $\Delta(\gamma) = f(y^*) - S_\gamma(y^*_\gamma)$ решения штрафных задач:

1. при $p \leq 1$ начиная с некоторого γ решения штрафных задач будет точно совпадать решениями исходной задачи ($\Delta(\gamma) = 0$);
2. при $p > 1$ для достаточно больших γ : $0 \leq \Delta(\gamma) \leq T / \gamma^{1/(p-1)}$, где $T = (L/\alpha)^{p/(p-1)} (1-1/p)/p^{1/(p-1)}$.

Доказательство можно найти в [6, 10].

Приведенные в теореме оценки подтверждают сделанные выше качественные выводы о характере влияния показателя степени p на характер и скорость сходимости метода штрафа.

Дополнительно рассмотрим вопрос о характере приближения оценок $y^*_{\gamma_k}$ к решению задачи при увеличении γ_k . Характер приближения изучим не в пространстве переменных, а по значениям функций $f(y^*_{\gamma_k})$ и $H(y^*_{\gamma_k})$.

Теорема 4.10. 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: $f(y^*_{\gamma_{k+1}}) \geq f(y^*_{\gamma_k})$, $H(y^*_{\gamma_{k+1}}) \leq H(y^*_{\gamma_k})$.

ДОКАЗАТЕЛЬСТВО. Рассмотрим два значения коэффициента штрафа $\gamma_{k+1} \geq \gamma_k$. Тогда будут справедливы два неравенства

$$f(y^*_{\gamma_k}) + \gamma_k H(y^*_{\gamma_k}) \leq f(y^*_{\gamma_{k+1}}) + \gamma_k H(y^*_{\gamma_{k+1}}),$$

$$f(y^*_{\gamma_{k+1}}) + \gamma_{k+1} H(y^*_{\gamma_{k+1}}) \leq f(y^*_{\gamma_k}) + \gamma_{k+1} H(y^*_{\gamma_k}).$$

После преобразования каждого из них получим, что

$$0 \leq (f(y^*_{\gamma_{k+1}}) - f(y^*_{\gamma_k})) + \gamma_k (H(y^*_{\gamma_{k+1}}) - H(y^*_{\gamma_k})),$$

$$(f(y^*_{\gamma_{k+1}}) - f(y^*_{\gamma_k})) + \gamma_{k+1} (H(y^*_{\gamma_{k+1}}) - H(y^*_{\gamma_k})) \leq 0.$$

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 Q и не является точкой безусловного минимума целевой функции, из теоремы следует, что решения вспомогательных штрафных задач приближаются к точке решения извне этой области.

4.7.3. Недостаточность локальных методов при использовании метода штрафов

В теореме об условиях сходимости метода штрафов предполагается, что при решении каждой задачи со штрафом определяется оценка $y_{\gamma_k}^*(\varepsilon_k)$ глобального минимума этой 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 $y_{\gamma_{k+1}}^*(\varepsilon_{k+1})$ 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 Q . 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, изменение значения коэффициента штрафа никак не повлияет на наличие этого локального минимума у функции штрафа. Следовательно, при достаточно большом коэффициенте γ у $S_\gamma(y)$ — функции штрафной задачи, образуется локальный минимум, расположенный в малой окрестности локального минимума функции штрафа. Если такая ситуация возникла на первой же итерации метода штрафов, то при последующем использовании локальных методов оценка решения останется в окрестности локального минимума и не будет приближаться к точному решению. Характерным признаком, по которому вычислительный метод может распознать эту ситуацию, является неограниченный рост значений функции $S_\gamma(y)$, вычисляемых в точках получаемых оценок, при $\gamma \rightarrow \infty$.

На рис.4.21 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 $(y_1^2 + y_2^2 - 11)^2 + (y_1 + y_2 - 7)^2 + (y_1 - y_2) \leq 0.3$. 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 Хука–Дживса, 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 — значению $\gamma = 1000$.

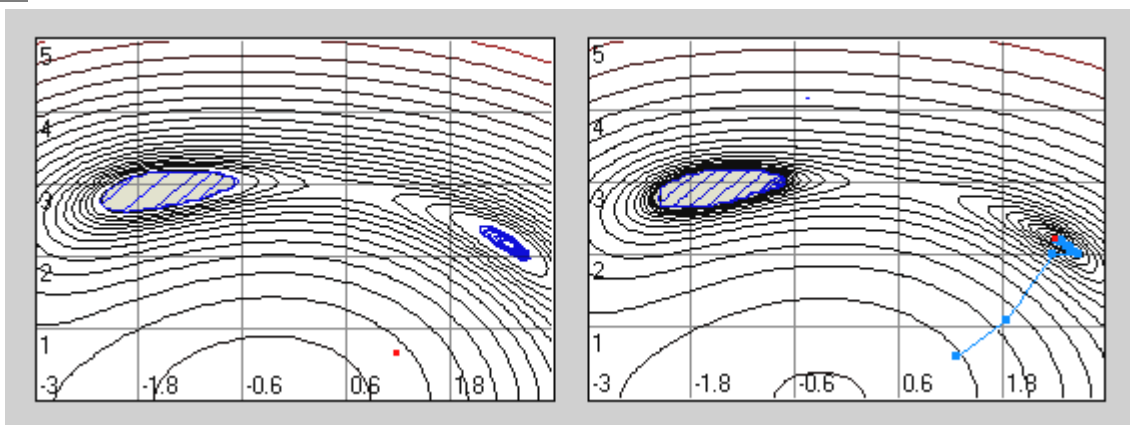


Рис. 4. 21. Возможность потери решения при наличии локального минимума у функции штрафа

4.7.4. Сочетание локальных методов с методами покрытий области

Большинство прикладных задач, связанных с поиском оптимальных решений, включают ограничения на переменные, заданные в виде неравенств достаточно общего вида. Одним из общих методов учета таких ограничений является метод внешнего штрафа, рассмотренный в разделе 4.7.1. Его применение, как это следует из теоремы 4.6 и обсуждения в разделе 4.7.3, требует применения методов многоэкстремальной оптимизации, осуществляющих, в общем случае, покрытие области поиска точками измерений для получения оценки глобального экстремума. В простейшем случае можно построить равномерное покрытие, применив метод Монте–Карло или регулярную сетку, например, ЛП_r–покрытие [18]. Однако для достижения необходимой точности при этом потребуется значительное количество измерений. Если функции задачи вычисляются достаточно быстро, то такой вариант организации решения будет возможен. Если же вычисление функций задачи дорогостояще, требуется применение других методов, которые способны строить адаптивные покрытия областей поиска точками, размещаемыми с различной плотностью в разных подобластях (более плотно в окрестности решений). Такие методы были рассмотрены в главе 3. При этом целесообразно получить лишь предварительную оценку решения, а для его уточнения использовать процедуры локальной оптимизации.

В том случае, когда размерность задачи будет слишком высока для использования методов, строящих адаптивные покрытия, единственным способом ее решения останутся методы локальной оптимизации, сочетаемые с простыми методами покрытия области для выбора начальных оценок решения.

Краткий обзор главы

В главе рассмотрен широкий круг вопросов по вычислительным методам поиска локально–оптимальных решений, а также математическим моделям, лежащим в основе их построения. Приведена классификация методов. Рассмотрены варианты всех основных групп методов для задач без ограничений (методы второго, первого порядка, методы прямого поиска), специальные версии этих методов с учетом двусторонних ограничений на переменные, общие методы учета функциональных ограничений–неравенств на примере метода штрафов. Приведен теоретический материал по рассмотренному кругу вопросов, включающий несколько теорем и лемм в которых анализируются основные свойства методов, приведены описания алгоритмов. В главе также содержится неформальное обсуждение свойств рассматриваемых методов, приведены примеры их использования, отмечены связанные с ними важные вычислительные эффекты. В тексте главы приведено большое количество иллюстраций, облегчающих восприятие материала. Часть иллюстраций получена с

использованием созданной программной лаборатории по методам локальной оптимизации «LocOpt». Необходимая практическая отработка обширного теоретического материала, приведенного в главе, может полностью выполняться с использованием этой программной лаборатории по разработанной программе лабораторного практикума.



Контрольные вопросы и упражнения

1. В чем заключается проблема поиска локального минимума? Дайте определение локального минимума, многоэкстремальной и одноэкстремальной задач.
2. Приведите описание общей структуры методов локальной оптимизации. Поясните понятия основного и рабочего шага.
3. Какая локальная информация может измеряться в точках испытаний? Приведите классификацию методов локального поиска с точки зрения измеряемой информации. Какова роль априорной модели задачи в интерпретации измерений?
4. Опишите постановку задачи определения величины перемещения вдоль выбранного методом направления. Приведите основные этапы вычислительного алгоритма для определения величины одномерных перемещений.
5. Приведите основные расчетные формулы для классических алгоритмов локальной оптимизации в задачах без ограничений: метода наискорейшего градиентного поиска и метода Ньютона. Что Вы можете сказать о свойствах этих методов?
6. Как влияет регулировка величины шага на свойства метода Ньютона?
7. Почему возникает задача коррекции матриц вторых производных и их оценок в методах локальной оптимизации? Опишите основную идею модифицированного алгоритма Холесского для коррекции матриц.
8. Приведите один из алгоритмов построения оценок матриц вторых производных в квазиньютоновских методах локальной оптимизации. Что такое квазиньютоновское условие?
9. Дайте сравнительное описание методов переменной метрики (квазиньютоновских методов) и метода растяжения пространства Шора Н.З.
10. Какие направления называют сопряженными? Приведите их основные свойства.
11. В чем заключается метод сопряженных градиентов Флетчера–Ривса? Сравните его свойства со свойствами метода наискорейшего градиентного поиска.
12. Какие методы относятся к методам прямого поиска? Приведите примеры таких методов.
13. Объясните различия между общими и специальными методами учета ограничений.
14. В чем заключаются принципы и особенности специального учета двусторонних ограничений на переменные в методах гладкой оптимизации? Поясните это на примере метода сопряженных градиентов и квазиньютоновском методе.
15. Как можно решить задачу с дополнительными ограничениями–неравенствами, используя метод внешнего штрафа? Как управлять гладкостью штрафа, как она может влиять на работу численных методов и на скорость сходимости метода штрафов?
16. Какие факторы могут привести к потере решения при использовании метода внешнего штрафа в сочетании с локальной оптимизацией?
17. Какие методы локальной оптимизации кажутся Вам наиболее перспективными в задачах поиска локально–оптимальных решений?

Предметный указатель

Алгоритм

- метода внешнего штрафа, 104
- метода Нелдера–Мида, 95
- метода переменной метрики, 87
 - модифицированный, 87
- метода растяжения пространства, 90
- метода Флетчера–Ривса, 93
- метода Хука–Дживса, 97
- Ньютона–Равсона с модификацией матрицы, 83
- одномерного поиска, 70
- определения промежутка, 70

Антиградиент, 68

Задача

- без ограничений, 66
- локальной оптимизации, 66
- математического программирования, 65
- одноэкстремальная, 66
- с ограничениями, 65
- со штрафом, 102

Задача оптимизации, классификация, 9

- многокритериальная, 9
- нелинейного программирования, 9

Задача оптимизации, классификация,

- безусловная, 48
- математического программирования, 25, 26
- нелинейного программирования, 48
- общая постановка, 25
- униmodalная, 31

Задача оптимизации, модель, 5

- вектор варьируемых параметров, 5, 6
- вектор допусков, 7
- критерии эффективности, 7
- область допустимых решений, 5, 7
- область поиска, 6
- ограничения, 5, 7

Задача оптимизации, примеры, 16

- определение координат манипулятора, 18
- определение местоположения и параметров движения объекта, 16

Задачи с ограничениями,

- индексная схема, 44

Испытания, шаги

- основные, 66
- рабочие, 66

Коэффициент штрафа, 102

Критерии

- выбора одномерного шага, 69

Метод

- к-го порядка, 68
- внешнего штрафа, 102

- второго порядка, 68
- градиентный, 71
- золотого сечения, 70
- квазиньютоновский, 70, 84
- наискорейшего градиентного поиска, 72
- Нелдера–Мида, 95
- Ньютона, 74
- Ньютона–Равсона, 77
- первого порядка, 68
- переменной метрики, 84
- покрытий, 108
- прямого поиска, 68
- растяжения пространства, 88
- сопряженных градиентов Флетчера–Ривса, 92
- сопряженных направлений, 91
- учета двусторонних ограничений на переменные, 99
- Хука–Дживса, 96
- Методы локальной оптимизации, обзор, 9
 - безусловная оптимизация, 9
 - условная оптимизация, 9
- Методы многомерной оптимизации, обзор, 11
 - адаптивное разбиение области поиска, 11
 - редукция размерности, 11
- Методы многоэкстремальной оптимизации, обзор, 10
 - адаптивные покрытия, 10
 - интервальные алгоритмы, 10
 - информационно-статистические алгоритмы, 10
 - равномерные покрытия, 10
- Методы оптимизации, оптимальность,
 - асимптотически оптимальные, 35
 - одношагово-оптимальные, 34, 38
 - оптимальные, 32
 - последовательно-оптимальные, 33
- Методы оптимизации, параллельные, 11
- Методы оптимизации, структура
 - оценка минимума, 6
 - решающее правило, 6
 - условие остановки, 6
 - характеристическая представимость, 11
- Методы оптимизации, структура,
 - оценка минимума, 28
 - решающее правило, 28
 - условие остановки, 28, 40
 - формальная модель, 28
 - характеристическая представимость, 35
- Минимум
 - глобальный, 65
 - локальный, 65
- Многошаговая схема редукции,
 - выпуклые ограничения, 59
 - липшицевость, 62

- монотонно-унимодальные ограничения, 61
- области с вычислимой границей, 58
- основное соотношение, 54
- сепарабельный случай, 62
- структура одномерных областей, 58
- Модель
 - задачи, 67
 - квадратичная, 76
- Направления
 - сопряженные, 90
 - к линейному многообразию, 90
- Ограничения
 - двусторонние, 98
 - общего вида, 101
 - специальные, 101
- Принцип локального спуска, 66
- Процесс оптимизации, основные понятия
 - испытание, 5
 - итерация, 6, 13, 15
 - минимизирующая последовательность, 6, 13
 - плотность итераций, 6
- Процесс оптимизации, основные понятия,
 - испытание, 27
 - минимизирующая последовательность, 30
- Разложение матрицы
 - спектральное, 78
 - Холесского, 81
 - модифицированное, 82
- Редукция,
 - кривые Пеано, 50
 - многошаговая схема, 54
 - размерности, 50
 - сложности, 49
- Решение
 - локально–оптимальное, 65
- Система имитации, 12, 13
- Сходимость
 - квадратичная, 75
 - линейная, 75
 - сверхлинейная, 75
- Сходимость,
 - всюду плотная, 41
 - глобальная, 44
 - двусторонняя, 38
 - локально-оптимальная, 42
 - метода оптимизации, 30
- Условие
 - квазиньютоново, 84
- Учебно-исследовательская система, 12
 - объект обучения, 12
 - объект показа, 12

сценарий обучения, 12

Формула

Бройдена, 85

Бройдена–Флетчера–Гольдфарба–Шанно, 85

Девидона–Флетчера–Пауэлла, 85

Функция

выпуклая (вниз), 67

квадратичная, 73

липшицева, 67

униmodalная, 70

штрафа внешнего, 102

Характеристические методы оптимизации,

информационно-статистический алгоритм глобального поиска (АГП), 37, 41

метод ломаных (Пиявского), 37, 40

метод последовательного сканирования (перебора), 36, 40

определение, 35

Литература

Литература к главе 1

1. Моисеев Н.Н. Математические методы системного анализа. М.:Наука,1981.
2. Гермейер Ю.Б. Введение в теорию исследования операций. М.: Наука, 1971.
3. Сухарев А.Г. Оптимальный поиск экстремума. М.: Изд-во МГУ, 1975.
4. Черноусько Ф.Л., Меликян А.А. Игровые задачи управления и поиска. М.: Наука, 1978.
5. Уайлд Д.Дж. Методы поиска экстремума. М.: Наука, 1967.
6. Моисеев Н.Н., Иванюков Ю.П., Столярова Е.М.. Методы оптимизации. М.: Наука, 1978.
7. Пшеничный Б.Н., Данилин Ю.М. Численные методы в экстремальных задачах. М.: Наука, 1975.
8. Химмельблау Д. Прикладное нелинейное программирование. М.: Мир, 1975.
9. Карманов В.Г. Математическое программирование. М.: Наука, 1986.
10. Малков В.П., Угодчиков А.Г. Оптимизация упругих систем. М.: Наука.1981.
11. Гилл Ф., Мюррей У., Райт М. Практическая оптимизация. М.: Мир, 1985.
12. Васильев Ф.П. Численные методы решения экстремальных задач. М.: Наука, 1980.
13. Фиакко А., Мак-Кормик Г. Нелинейное программирование. Методы последовательной безусловной оптимизации. М.: Мир, 1972.
14. Базара М., Шетти К. Нелинейное программирование. Теория и алгоритмы. М.: Мир, 1982.
15. Численные методы условной оптимизации. / Под ред. Ф.Гилла и У.Мюррея. М.: Мир,1977.
16. Евтушенко Ю.Г. Методы решения экстремальных задач и их применение в системах оптимизации. М.: Наука, 1982.
17. Бертсекас Д. Условная оптимизация и методы множителей Лагранжа. М.: Радио и связь, 1987.
18. Сухарев А.Г., Тимофеев А.В., Федоров В.В. Курс методов оптимизации. М.: Наука, 1986. 325с.
19. Тихонов А.Н., Арсенин В.Я. Методы решения некорректных задач. М.: Наука, 1979.
20. Батищев Д.И. Поисковые методы оптимального проектирования. М.: Советское радио, 1975.
21. Растринин Л.А. Статистические методы поиска. Наука,1968.
22. Соболев И.М., Статников Р.Б. Выбор оптимальных параметров в задачах со многими критериями. М.: Наука, 1981. 110с.
23. Жиглявский А.А. Математическая теория глобального случайного поиска. М.: Изд-во ЛГУ, 1985. 293 с.
24. Чичинадзе В.К. Решение невыпуклых нелинейных задач оптимизации. М.: 1983, 256 с.
25. Стронгин Р.Г. Численные методы в экстремальных задачах. М.: Наука, 1978.
26. Стронгин Р.Г. Поиск глобального оптимума. М.: Знание, 1990.
27. Жиглявский А.А., Жилинскас А.Г. Методы поиска глобального экстремума. М.: Наука, 1991.
28. Strongin R.G., Sergeyev Ya.D. Global Optimization with Non-Convex Constraints. Sequential and Parallel Algorithms. –Dordrecht: Kluwer Academic Publishers. The Netherlands, 2000. 728 pp.

29. Horst R., Tuy H. Global Optimization – Deterministic Approaches. –Berlin: Springer, 1990.
30. Horst R., Pardalos P.M., ets. Handbook of Global Optimization. –Dordrecht: Kluwer, 1995.
31. Pinter J. Global optimization in Action. –Dordrecht: Kluwer, 1996.
32. Кони Р., Райфа Х. Принятие решений при многих критериях: предпочтения и замечания. – М: Радио и связь, 1981.
33. Подиновский В.В., Ногин В.Д. Парето-оптимальные решения многокритериальных задач. М.: Наука, 1982. 254 с.
34. Альфред Г., Хорцвергер Ю. Введение в интервальные вычисления. М.: Мир, 1987.
35. Grishagin V.A., Sergeyev Ya.D., Strongin R.G. Parallel characteristical algorithms for solving problems of global optimization//Jurnal of global optimization. 1997, v.10, p.185-206.
36. Strongin R.G., Markin D.L., Markina M.V. Reduction of Multi-extremum Multi-criterion Problems with Constraints to Unconstrained Optimization Problems: Theory and Algorithms. Computational and Mathematical Modelling, V. 6. № 4. 1995, P.242-248. Plenum Publishing Corporation.
37. Гергель В.П., Стронгин Р.Г. АБСОЛЮТ. Программная система для исследования и изучения методов глобальной оптимизации. Учебное пособие. Нижний Новгород: Изд-во Нижегородского университета. 1998.
38. Малков В.П., Маркина М.В. Поэтапная параметрическая оптимизация. Учебное пособие. Н. Новгород. Изд-во Нижегородского университета. 1998.
39. Белоглазов И.Н., Джанджава Г.И., Чичин Г.П. Основы навигации по геофизическим полям /Под ред.А.А.Красовского.
40. Цыпкин Я.З., Поляк Б.Т. Огрубленный метод максимального правдоподобия. Динамика систем. Математические методы теории колебаний. Межвузовский сборник. Вып.12. Горький. Изд-во ГГУ. 1977. стр.22-46.

Литература к главе 2

1. Базара М., Шетти К. Нелинейное программирование. Теория и алгоритмы.-М.:Мир, 1982.
2. Васильев Ф.П. Численные методы решения экстремальных задач.-М.:Наука, 1980.
3. Сухарев А.Г. Оптимальный поиск экстремума. – М.: Изд МГУ, 1975.
4. Стронгин Р.Г. Численные методы в многоэкстремальных задачах. Информационно-статистический подход. М.: Наука, 1978.
5. Strongin R.G., Sergeyev Ya.D. Global optimization with non-convex constraints: Sequential and parallel algorithms, Kluwer Academic Publishers, Dordrecht, Netherlands, 2000.
6. Пиявский С.А. Один алгоритм отыскания абсолютного минимума функции.-ЖВМ и МФ, №4, 1972, с.888-896.
7. Батищев Д.И. Поисковые методы оптимального проектирования.-М.: Советское радио,1975.
8. Гермейер Ю.Б. Введение в теорию исследования операций. - М.:Наука,1971.
9. Сухарев А.Г. Наилучшие стратегии последовательного поиска экстремума.-ЖВМ и МФ, т.12, №1, 1972, с.35-50.
10. Kiefer J. Sequential minimax search for a maximum.-Proc. Amer. Math. Soc., Vol.4, No.3, 1953, 502-506.
11. Сухарев А.Г. Глобальный экстремум и методы его отыскания.-В кн.: Математические методы в исследовании операций.-М.: Изд.МГУ, 1981, с.4-37.
12. Уайлд Д.Дж. Методы поиска экстремума.-М.:Наука, 1967.

13. Евтушенко Ю.Г. Методы решения экстремальных задач и их применение в системах оптимизации. - М.:Наука, 1982.
14. Иванов В.В. Об оптимальных алгоритмах минимизации функций некоторых классов.- Кибернетика, 1972, №4, с.81-94.
15. Леонов В.В. Метод покрытий для отыскания глобального максимума функций от многих переменных.-В кн.: исследования по кибернетике. По ред.Ляпунова А.А.- М.: Советское радио, 1970, с.41-52.
16. Ганшин Г.С. Оптимальные пассивные алгоритмы вычисления наибольшего значения функций на отрезке.-ЖВМ и МФ, 1977, т.17, №3, с.562-571.
17. Гирлин С.К. Об оптимальных по точности интерполяции и минимизации функций классов $C_{2,L_1,\dots,L_m,N}$.- Известия вузов. Математика, 1978, №10, с.95-98.
18. Зализняк И.Ф., Лигун А.А. Об оптимальных стратегиях поиска глобального максимума функции. – ЖВМ и МФ, 1978,т.18, №2, с.314-321.
19. Певный А.Б. Об оптимальных стратегиях поиска максимума функции с ограниченной старшей производной.- ЖВМ и МФ, 1982, т.22, №5, с.1061-1066.
20. Тарасова В.П. Оптимальные стратегии поиска области наибольших значений для некоторого класса функций. – ЖВМ и МФ, 1978, т.18, №4, с.886-896.
21. Моцкус Й.Б. О байесовых методах поиска экстремума. - Автоматика и вычислительная техника, 1972, № 3, с.53-62.
22. Жилинскас А.Г. , Моцкус Й.Б. Об одном байесовом методе поиска минимума.- Автоматика и вычислительная техника, 1972, № 4, с.42-44.
23. Жилинскас А.Г. , Моцкус Й.Б., Тимофеев Л.Л. Байесов метод поиска экстремума с ограниченной памятью. - Автоматика и вычислительная техника, 1972, № 6, с.37-42.
24. Черноусько Ф.Л., Меликян А.А. Игровые задачи управления и поиска.- Мю:Наука, 1978.
25. Коротченко А.Г. Об одном алгоритме поиска наибольшего значения одномерных функций. – ЖВМ и МФ, 1978, т.18, № 3, с.563-573.
26. Kushner H. A New Method of Locating the Maximum Point of an Arbitrary Multippeak Curve in the Presence of Noise. – Transactions ASME, Ser. D, J. Basic Eng., 1964, vol.86, No.1, p.97-106.
27. Шалтянис В.Р. Об одном методе многоэкстремальной оптимизации. - Автоматика и вычислительная техника, 1971, № 3, с.33-38.
28. Жилинскас А.Г. Одношаговый байесовский метод поиска экстремума функций одной переменной.- Кибернетика, 1975, № 1, с.139-144.
29. Городецкий С.Ю., Неймарк Ю.И. О поисковых характеристиках алгоритма глобальной оптимизации с адаптивной стохастической моделью.- В кн.: Проблемы случайного поиска, вып. 9 –Рига: Изд. Зинатне, 1981, с.83-105.
30. Törn A.A., Žilinskas A. Global Optimization. Lecture Notes in Computer Science 350. Springer Verlag, Berlin, 1989.
31. Гришагин В.А. Об условиях сходимости для одного класса алгоритмов глобального поиска. – В кн.: Тез. докл. III Всес. семинара "Численные методы нелинейного программирования". –Харьков: Изд ХГУ, 1979, с. 82-84.
32. Pinter J. Global Optimization in Action (Continuous and Lipschitz Optimization: Algorithms, Implementations and Applications). Kluwer Academic Publishers, Dordrecht, 1996.
33. 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.
34. Sergeyev Ya.D. On Convergence of "Divide the Best" Global Optimization Algorithms. – Optimization, Vol.44, 1998, pp.303-325.

35. Sergeyev Ya.D. An information global optimization algorithm with local tuning. SIAM Journal on Optimization, 5, 4, 1995, 17-31.
36. Gergel V.P., Sergeyev Ya.D. Sequential and parallel global optimization methods using derivatives. Computers & Mathematics with Applications, 37, 4/5, 163-180.
37. Strongin R.G., Markin D.L. Minimization of multiextremal functions with nonconvex constraints. Cybernetics 22(4), 1986, p.486-493.
38. Стронгин Р.Г. Поиск глобального оптимума. Серия "Математика и кибернетика" 2, М.:Знание, 1990.

Литература к главе 3.

1. Ермаков С.М. Метод Монте-Карло и смежные вопросы. – М.:Наука, 1971.
2. Соболев И.М. Численные методы Монте-Карло. – М.:Наука, 1973.
3. Соболев И.М., Статников Р.Б. ЛП-поиск и задачи оптимального конструирования. – В кн.: Проблемы случайного поиска. Вып. 1.- Рига: изд. Зинатне, 1973, с.117-135.
4. Сухарев А.Г. Глобальный экстремум и методы его отыскания.-В кн.: Математические методы в исследовании операций.-М.: Изд.МГУ, 1981, с.4-37.
5. Карр Ч., Хоув Ч. Количественные методы принятия решений в управлении и экономике. – М.:Мир, 1966.
6. Стронгин Р.Г. (1978) Численные методы в многоэкстремальных задачах. Информационно- статистический подход. М.: Наука.
7. Стронгин Р.Г. Поиск глобального оптимума. Серия "Математика и кибернетика" 2, М.:Знание, 1990.
8. Strongin R.G., Sergeyev Ya.D. Global optimization with non-convex constraints: Sequential and parallel algorithms, Kluwer Academic Publishers, Dordrecht, Netherlands, 2000.
9. Sergeyev Ya.D., Grishagin V.A. Parallel asynchronous global search and the nested optimization scheme, Journal of Computational Analysis & Applications, 3(2), 2001, pp.123-145.
10. Gergel V.P., Strongin R.G. Multiple Peano Curves in Recognition Problems. Pattern Recognition and Image Analysis, 2(2), 1992, 161-164.
11. Карманов В.Г. Математическое программирование. – М.:Наука, 1975.
12. Стронгин Р.Г., Гришагин В.А. Оптимизация многоэкстремальных функций при монотонно унимодальных ограничениях. - Изв. АН СССР. Техническая кибернетика, 1984, №4.
13. Strongin R.G., Markin D.L. Minimization of multiextremal functions with nonconvex constraints. Cybernetics 22(4), 1986, p.486-493.
14. Horst R., Pardalos P.M. Handbook on Global Optimization. Kluwer Academic Publishers, Dordrecht, 1995.

Литература к главе 4

1. Пшеничный Б.Н., Данилин Ю.М. Численные методы в экстремальных задачах. –М.: Наука, 1975. –319 с.
2. Гилл Ф., Мюррей У., Райт М. Практическая оптимизация. –М.: Мир, 1985. – 509 с.
3. Уайлд Д.Дж. Методы поиска экстремума. –М.: Наука, 1967. –267 с.
4. Габасов Р. Кириллова Ф. Методы оптимизации. –Минск: Изд-во БГУ, 1975.
5. Черноусько Ф.Л., Меликян А.А. Игровые задачи управления и поиска. –М.: Наука, 1978.

6. Сухарев А.Г., Тимофеев А.В., Федоров В.В. Курс методов оптимизации. –М.: Наука, 1986. –325с.
7. Моисеев Н.Н., Иванилов Ю.П., Столярова Е.М.. Методы оптимизации. –М.: Наука, 1978. –350 с.
8. Аоки М. Введение в оптимизацию. — М.: Наука, 1975. –343 с.
9. Васильев Ф.П. Численные методы решения экстремальных задач. –М.: Наука, 1980. –518 с.
10. Карманов В.Г. Математическое программирование. –М.: Наука, 1986. –285 с.
11. Базара М., Шетти К. Нелинейное программирование. Теория и алгоритмы. –М.: Мир, 1982. –583 с.
12. Химмельблау Д. Прикладное нелинейное программирование. –М.: Мир, 1975.
13. Банди Б. Методы оптимизации. Вводный курс. М.: Радио и связь. 1988. –128 с.
14. Фиакко А., Мак-Кормик Г. Нелинейное программирование. Методы последовательной безусловной оптимизации. –М.: Мир, 1972.
15. Бертсекас Д. Условная оптимизация и методы множителей Лагранжа. –М.: Радио и связь, 1987. –399 с.
16. Евтушенко Ю.Г. Методы решения экстремальных задач и их применение в системах оптимизации. –М.: Наука, 1982. –432 с.
17. Численные методы условной оптимизации. / Под ред. Ф.Гилла и У.Мюррея. –М.: Мир, 1977. –290 с.
18. Соболев И.М., Статников Р.Б. Выбор оптимальных параметров в задачах со многими критериями. –М.: Наука, 1981. –110с.