

If a distribution has a strong central tendency, so that most of its area is under a single peak, then the median is an estimator of the central value. It is a more robust estimator than the mean is: The median fails as an estimator only if the area in the tails is large, while the mean fails if the first moment of the tails is large; it is easy to construct examples where the first moment of the tails is large even though their area is negligible.

To find the median of a set of values, one can proceed by sorting the set and then applying (14.1.14). This is a process of order  $N \log N$ . You might rightly think that this is wasteful, since it yields much more information than just the median (e.g., the upper and lower quartile points, the deciles, etc.). In fact, we saw in §8.5 that the element  $x_{(N+1)/2}$  can be located in of order  $N$  operations. Consult that section for routines.

The *mode* of a probability distribution function  $p(x)$  is the value of  $x$  where it takes on a maximum value. The mode is useful primarily when there is a single, sharp maximum, in which case it estimates the central value. Occasionally, a distribution will be *bimodal*, with two relative maxima; then one may wish to know the two modes individually. Note that, in such cases, the mean and median are not very useful, since they will give only a “compromise” value between the two peaks.

#### CITED REFERENCES AND FURTHER READING:

- Bevington, P.R. 1969, *Data Reduction and Error Analysis for the Physical Sciences* (New York: McGraw-Hill), Chapter 2.
- Stuart, A., and Ord, J.K. 1987, *Kendall's Advanced Theory of Statistics*, 5th ed. (London: Griffin and Co.) [previous eds. published as Kendall, M., and Stuart, A., *The Advanced Theory of Statistics*], vol. 1, §10.15
- Norusis, M.J. 1982, *SPSS Introductory Guide: Basic Statistics and Operations*; and 1985, *SPSS-X Advanced Statistics Guide* (New York: McGraw-Hill).
- Chan, T.F., Golub, G.H., and LeVeque, R.J. 1983, *American Statistician*, vol. 37, pp. 242–247. [1]
- Cramér, H. 1946, *Mathematical Methods of Statistics* (Princeton: Princeton University Press), §15.10. [2]

## 14.2 Do Two Distributions Have the Same Means or Variances?

Not uncommonly we want to know whether two distributions have the same mean. For example, a first set of measured values may have been gathered before some event, a second set after it. We want to know whether the event, a “treatment” or a “change in a control parameter,” made a difference.

Our first thought is to ask “how many standard deviations” one sample mean is from the other. That number may in fact be a useful thing to know. It does relate to the strength or “importance” of a difference of means *if that difference is genuine*. However, by itself, it says nothing about whether the difference *is* genuine, that is, statistically significant. A difference of means can be very small compared to the standard deviation, and yet very significant, if the number of data points is large. Conversely, a difference may be moderately large but not significant, if the data

are sparse. We will be meeting these distinct concepts of *strength* and *significance* several times in the next few sections.

A quantity that measures the significance of a difference of means is not the number of standard deviations that they are apart, but the number of so-called *standard errors* that they are apart. The standard error of a set of values measures the accuracy with which the sample mean estimates the population (or “true”) mean. Typically the standard error is equal to the sample’s standard deviation divided by the square root of the number of points in the sample.

### Student’s *t*-test for Significantly Different Means

Applying the concept of standard error, the conventional statistic for measuring the significance of a difference of means is termed *Student’s t*. When the two distributions are thought to have the same variance, but possibly different means, then Student’s *t* is computed as follows: First, estimate the standard error of the difference of the means,  $s_D$ , from the “pooled variance” by the formula

$$s_D = \sqrt{\frac{\sum_{i \in A} (x_i - \bar{x}_A)^2 + \sum_{i \in B} (x_i - \bar{x}_B)^2}{N_A + N_B - 2} \left( \frac{1}{N_A} + \frac{1}{N_B} \right)} \quad (14.2.1)$$

where each sum is over the points in one sample, the first or second, each mean likewise refers to one sample or the other, and  $N_A$  and  $N_B$  are the numbers of points in the first and second samples, respectively. Second, compute  $t$  by

$$t = \frac{\bar{x}_A - \bar{x}_B}{s_D} \quad (14.2.2)$$

Third, evaluate the significance of this value of  $t$  for Student’s distribution with  $N_A + N_B - 2$  degrees of freedom, by equations (6.4.7) and (6.4.9), and by the routine `betai` (incomplete beta function) of §6.4.

The significance is a number between zero and one, and is the probability that  $|t|$  could be this large or larger just by chance, for distributions with equal means. Therefore, a small numerical value of the significance (0.05 or 0.01) means that the observed difference is “very significant.” The function  $A(t|\nu)$  in equation (6.4.7) is one minus the significance.

As a routine, we have

```

SUBROUTINE tttest(data1,n1,data2,n2,t,prob)
  INTEGER n1,n2
  REAL prob,t,data1(n1),data2(n2)
C  USES avevar,betai
      Given the arrays data1(1:n1) and data2(1:n2), this routine returns Student's t as t,
      and its significance as prob, small values of prob indicating that the arrays have significantly
      different means. The data arrays are assumed to be drawn from populations with the same
      true variance.
  REAL ave1,ave2,df,var,var1,var2,betai
  call avevar(data1,n1,ave1,var1)
  call avevar(data2,n2,ave2,var2)
  df=n1+n2-2
  var=((n1-1)*var1+(n2-1)*var2)/df
  t=(ave1-ave2)/sqrt(var*(1./n1+1./n2))
  prob=betai(0.5*df,0.5,df/(df+t**2))
  return
END

```

Degrees of freedom.  
Pooled variance.  
See equation (6.4.9).

which makes use of the following routine for computing the mean and variance of a set of numbers,

```

SUBROUTINE avevar(data,n,ave,var)
INTEGER n
REAL ave,var,data(n)
  Given array data(1:n), returns its mean as ave and its variance as var.
INTEGER j
REAL s,ep
ave=0.0
do 11 j=1,n
  ave=ave+data(j)
enddo 11
ave=ave/n
var=0.0
ep=0.0
do 12 j=1,n
  s=data(j)-ave
  ep=ep+s
  var=var+s*s
enddo 12
var=(var-ep**2/n)/(n-1)      Corrected two-pass formula (14.1.8).
return
END

```

The next case to consider is where the two distributions have significantly different variances, but we nevertheless want to know if their means are the same or different. (A treatment for baldness has caused some patients to *lose* all their hair and turned others into werewolves, but we want to know if it helps cure baldness *on the average!*) Be suspicious of the unequal-variance *t*-test: If two distributions have very different variances, then they may also be substantially different in shape; in that case, the difference of the means may not be a particularly useful thing to know.

To find out whether the two data sets have variances that are significantly different, you use the *F*-test, described later on in this section.

The relevant statistic for the unequal variance *t*-test is

$$t = \frac{\bar{x}_A - \bar{x}_B}{[\text{Var}(x_A)/N_A + \text{Var}(x_B)/N_B]^{1/2}} \quad (14.2.3)$$

This statistic is distributed *approximately* as Student's *t* with a number of degrees of freedom equal to

$$\frac{\left[ \frac{\text{Var}(x_A)}{N_A} + \frac{\text{Var}(x_B)}{N_B} \right]^2}{\frac{[\text{Var}(x_A)/N_A]^2}{N_A - 1} + \frac{[\text{Var}(x_B)/N_B]^2}{N_B - 1}} \quad (14.2.4)$$

Expression (14.2.4) is in general not an integer, but equation (6.4.7) doesn't care.

The routine is

```

SUBROUTINE tutest(data1,n1,data2,n2,t,prob)
INTEGER n1,n2
REAL prob,t,data1(n1),data2(n2)
C USES avevar,beta1
  Given the arrays data1(1:n1) and data2(1:n2), this routine returns Student's t as t,
  and its significance as prob, small values of prob indicating that the arrays have significantly

```

different means. The data arrays are allowed to be drawn from populations with unequal variances.

```

REAL ave1,ave2,df,var1,var2,betai
call avevar(data1,n1,ave1,var1)
call avevar(data2,n2,ave2,var2)
t=(ave1-ave2)/sqrt(var1/n1+var2/n2)
df=(var1/n1+var2/n2)**2/((var1/n1)**2/(n1-1)+(var2/n2)**2/(n2-1))
prob=betai(0.5*df,0.5,df/(df+t**2))
return
END

```

Our final example of a Student's  $t$  test is the case of *paired samples*. Here we imagine that much of the variance in *both* samples is due to effects that are point-by-point identical in the two samples. For example, we might have two job candidates who have each been rated by the same ten members of a hiring committee. We want to know if the means of the ten scores differ significantly. We first try `ttest` above, and obtain a value of `prob` that is not especially significant (e.g.,  $> 0.05$ ). But perhaps the significance is being washed out by the tendency of some committee members always to give high scores, others always to give low scores, which increases the apparent variance and thus decreases the significance of any difference in the means. We thus try the paired-sample formulas,

$$\text{Cov}(x_A, x_B) \equiv \frac{1}{N-1} \sum_{i=1}^N (x_{Ai} - \bar{x}_A)(x_{Bi} - \bar{x}_B) \quad (14.2.5)$$

$$s_D = \left[ \frac{\text{Var}(x_A) + \text{Var}(x_B) - 2\text{Cov}(x_A, x_B)}{N} \right]^{1/2} \quad (14.2.6)$$

$$t = \frac{\bar{x}_A - \bar{x}_B}{s_D} \quad (14.2.7)$$

where  $N$  is the number in each sample (number of pairs). Notice that it is important that a particular value of  $i$  label the corresponding points in each sample, that is, the ones that are paired. The significance of the  $t$  statistic in (14.2.7) is evaluated for  $N - 1$  degrees of freedom.

The routine is

```

SUBROUTINE tptest(data1,data2,n,t,prob)
INTEGER n
REAL prob,t,data1(n),data2(n)
C USES avevar,betai
    Given the paired arrays data1(1:n) and data2(1:n), this routine returns Student's t for
    paired data as t, and its significance as prob, small values of prob indicating a significant
    difference of means.
INTEGER j
REAL ave1,ave2,cov,df,sd,var1,var2,betai
call avevar(data1,n,ave1,var1)
call avevar(data2,n,ave2,var2)
cov=0.
do 11 j=1,n
    cov=cov+(data1(j)-ave1)*(data2(j)-ave2)
enddo 11
df=n-1
cov=cov/df
sd=sqrt((var1+var2-2.*cov)/n)
t=(ave1-ave2)/sd

```

```

prob=betai(0.5*df,0.5,df/(df+t**2))
return
END

```

## ***F-Test for Significantly Different Variances***

The *F-test* tests the hypothesis that two samples have different variances by trying to reject the null hypothesis that their variances are actually consistent. The statistic  $F$  is the ratio of one variance to the other, so values either  $\gg 1$  or  $\ll 1$  will indicate very significant differences. The distribution of  $F$  in the null case is given in equation (6.4.11), which is evaluated using the routine `betai`. In the most common case, we are willing to disprove the null hypothesis (of equal variances) by either very large or very small values of  $F$ , so the correct significance is *two-tailed*, the sum of two incomplete beta functions. It turns out, by equation (6.4.3), that the two tails are always equal; we need compute only one, and double it. Occasionally, when the null hypothesis is strongly viable, the identity of the two tails can become confused, giving an indicated probability greater than one. Changing the probability to two minus itself correctly exchanges the tails. These considerations and equation (6.4.3) give the routine

```

SUBROUTINE ftest(data1,n1,data2,n2,f,prob)
INTEGER n1,n2
REAL f,prob,data1(n1),data2(n2)
C USES avevar,betai
    Given the arrays data1(1:n1) and data2(1:n2), this routine returns the value of f, and
    its significance as prob. Small values of prob indicate that the two arrays have significantly
    different variances.
REAL ave1,ave2,df1,df2,var1,var2,betai
call avevar(data1,n1,ave1,var1)
call avevar(data2,n2,ave2,var2)
if(var1.gt.var2)then      Make F the ratio of the larger variance to the smaller one.
    f=var1/var2
    df1=n1-1
    df2=n2-1
else
    f=var2/var1
    df1=n2-1
    df2=n1-1
endif
prob=2.*betai(0.5*df2,0.5*df1,df2/(df2+df1*f))
if(prob.gt.1.)prob=2.-prob
return
END

```

### CITED REFERENCES AND FURTHER READING:

- von Mises, R. 1964, *Mathematical Theory of Probability and Statistics* (New York: Academic Press), Chapter IX(B).  
 Norusis, M.J. 1982, *SPSS Introductory Guide: Basic Statistics and Operations*, and 1985, *SPSS-X Advanced Statistics Guide* (New York: McGraw-Hill).

Sample page from NUMERICAL RECIPES IN FORTRAN 77: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43064-X)  
 Copyright (C) 1986-1992 by Cambridge University Press. Programs Copyright (C) 1986-1992 by Numerical Recipes Software.  
 Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books, diskettes, or CDROMs visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to [trade@cup.cam.ac.uk](mailto:trade@cup.cam.ac.uk) (outside North America).

### 14.3 Are Two Distributions Different?

Given two sets of data, we can generalize the questions asked in the previous section and ask the single question: Are the two sets drawn from the same distribution function, or from different distribution functions? Equivalently, in proper statistical language, “Can we disprove, to a certain required level of significance, the null hypothesis that two data sets are drawn from the same population distribution function?” Disproving the null hypothesis in effect proves that the data sets are from different distributions. Failing to disprove the null hypothesis, on the other hand, only shows that the data sets can be *consistent* with a single distribution function. One can never *prove* that two data sets come from a single distribution, since (e.g.) no practical amount of data can distinguish between two distributions which differ only by one part in  $10^{10}$ .

Proving that two distributions are different, or showing that they are consistent, is a task that comes up all the time in many areas of research: Are the visible stars distributed uniformly in the sky? (That is, is the distribution of stars as a function of declination — position in the sky — the same as the distribution of sky area as a function of declination?) Are educational patterns the same in Brooklyn as in the Bronx? (That is, are the distributions of people as a function of last-grade-attended the same?) Do two brands of fluorescent lights have the same distribution of burn-out times? Is the incidence of chicken pox the same for first-born, second-born, third-born children, etc.?

These four examples illustrate the four combinations arising from two different dichotomies: (1) The data are either continuous or binned. (2) Either we wish to compare one data set to a known distribution, or we wish to compare two equally unknown data sets. The data sets on fluorescent lights and on stars are continuous, since we can be given lists of individual burnout times or of stellar positions. The data sets on chicken pox and educational level are binned, since we are given tables of numbers of events in discrete categories: first-born, second-born, etc.; or 6th Grade, 7th Grade, etc. Stars and chicken pox, on the other hand, share the property that the null hypothesis is a known distribution (distribution of area in the sky, or incidence of chicken pox in the general population). Fluorescent lights and educational level involve the comparison of two equally unknown data sets (the two brands, or Brooklyn and the Bronx).

One can always turn continuous data into binned data, by grouping the events into specified ranges of the continuous variable(s): declinations between 0 and 10 degrees, 10 and 20, 20 and 30, etc. Binning involves a loss of information, however. Also, there is often considerable arbitrariness as to how the bins should be chosen. Along with many other investigators, we prefer to avoid unnecessary binning of data.

The accepted test for differences between binned distributions is the *chi-square test*. For continuous data as a function of a single variable, the most generally accepted test is the *Kolmogorov-Smirnov test*. We consider each in turn.

#### Chi-Square Test

Suppose that  $N_i$  is the number of events observed in the  $i$ th bin, and that  $n_i$  is the number expected according to some known distribution. Note that the  $N_i$ 's are

Sample page from NUMERICAL RECIPES IN FORTRAN 77: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43064-X)  
 Copyright (C) 1986-1992 by Cambridge University Press. Programs Copyright (C) 1986-1992 by Numerical Recipes Software.  
 Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books, diskettes, or CDROMs visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to [trade@cup.cam.ac.uk](mailto:trade@cup.cam.ac.uk) (outside North America).