

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

integers, while the n_i 's may not be. Then the chi-square statistic is

$$\chi^2 = \sum_i \frac{(N_i - n_i)^2}{n_i} \quad (14.3.1)$$

where the sum is over all bins. A large value of χ^2 indicates that the null hypothesis (that the N_i 's are drawn from the population represented by the n_i 's) is rather unlikely.

Any term j in (14.3.1) with $0 = n_j = N_j$ should be omitted from the sum. A term with $n_j = 0$, $N_j \neq 0$ gives an infinite χ^2 , as it should, since in this case the N_i 's cannot possibly be drawn from the n_i 's!

The *chi-square probability function* $Q(\chi^2|\nu)$ is an incomplete gamma function, and was already discussed in §6.2 (see equation 6.2.18). Strictly speaking $Q(\chi^2|\nu)$ is the probability that the sum of the squares of ν random *normal* variables of unit variance (and zero mean) will be greater than χ^2 . The terms in the sum (14.3.1) are not individually normal. However, if either the number of bins is large ($\gg 1$), or the number of events in each bin is large ($\gg 1$), then the chi-square probability function is a good approximation to the distribution of (14.3.1) in the case of the null hypothesis. Its use to estimate the significance of the chi-square test is standard.

The appropriate value of ν , the number of degrees of freedom, bears some additional discussion. If the data are collected with the model n_i 's fixed — that is, not later renormalized to fit the total observed number of events ΣN_i — then ν equals the number of bins N_B . (Note that this is *not* the total number of *events*!) Much more commonly, the n_i 's are normalized after the fact so that their sum equals the sum of the N_i 's. In this case the correct value for ν is $N_B - 1$, and the model is said to have one constraint (knstrn=1 in the program below). If the model that gives the n_i 's has additional free parameters that were adjusted after the fact to agree with the data, then each of these additional “fitted” parameters decreases ν (and increases knstrn) by one additional unit.

We have, then, the following program:

```
void chsone(float bins[], float ebins[], int nbins, int knstrn, float *df,
float *chsq, float *prob)
Given the array bins[1..nbins] containing the observed numbers of events, and an array
ebins[1..nbins] containing the expected numbers of events, and given the number of con-
straints knstrn (normally one), this routine returns (trivially) the number of degrees of freedom
df, and (nontrivially) the chi-square chsq and the significance prob. A small value of prob
indicates a significant difference between the distributions bins and ebins. Note that bins
and ebins are both float arrays, although bins will normally contain integer values.
{
float gammq(float a, float x);
void nrerror(char error_text[]);
int j;
float temp;

*df=nbins-knstrn;
*chsq=0.0;
for (j=1;j<=nbins;j++) {
if (ebins[j] <= 0.0) nrerror("Bad expected number in chsone");
temp=bins[j]-ebins[j];
*chsq += temp*temp/ebins[j];
}
*prob=gammq(0.5*(*df),0.5*(*chsq));      Chi-square probability function. See §6.2.
}
```

Next we consider the case of comparing *two* binned data sets. Let R_i be the number of events in bin i for the first data set, S_i the number of events in the same bin i for the second data set. Then the chi-square statistic is

$$\chi^2 = \sum_i \frac{(R_i - S_i)^2}{R_i + S_i} \quad (14.3.2)$$

Comparing (14.3.2) to (14.3.1), you should note that the denominator of (14.3.2) is *not* just the average of R_i and S_i (which would be an estimator of n_i in 14.3.1). Rather, it is twice the average, the sum. The reason is that each term in a chi-square sum is supposed to approximate the square of a normally distributed quantity with unit variance. The variance of the difference of two normal quantities is the sum of their individual variances, not the average.

If the data were collected in such a way that the sum of the R_i 's is necessarily equal to the sum of S_i 's, then the number of degrees of freedom is equal to one less than the number of bins, $N_B - 1$ (that is, $\text{knstrn} = 1$), the usual case. If this requirement were absent, then the number of degrees of freedom would be N_B . Example: A birdwatcher wants to know whether the distribution of sighted birds as a function of species is the same this year as last. Each bin corresponds to one species. If the birdwatcher takes his data to be the first 1000 birds that he saw in each year, then the number of degrees of freedom is $N_B - 1$. If he takes his data to be all the birds he saw on a random sample of days, the same days in each year, then the number of degrees of freedom is N_B ($\text{knstrn} = 0$). In this latter case, note that he is also testing whether the birds were more numerous overall in one year or the other: That is the extra degree of freedom. Of course, any additional constraints on the data set lower the number of degrees of freedom (i.e., increase knstrn to *more positive* values) in accordance with their number.

The program is

```
void chstwo(float bins1[], float bins2[], int nbins, int knstrn, float *df,
           float *chsq, float *prob)
Given the arrays bins1[1..nbins] and bins2[1..nbins], containing two sets of binned
data, and given the number of constraints knstrn (normally 1 or 0), this routine returns the
number of degrees of freedom df, the chi-square chsq, and the significance prob. A small value
of prob indicates a significant difference between the distributions bins1 and bins2. Note that
bins1 and bins2 are both float arrays, although they will normally contain integer values.
{
    float gammq(float a, float x);
    int j;
    float temp;

    *df=nbins-knstrn;
    *chsq=0.0;
    for (j=1;j<=nbins;j++)
        if (bins1[j] == 0.0 && bins2[j] == 0.0)
            --(*df);                                     No data means one less degree of free-
        else {
            temp=bins1[j]-bins2[j];
            *chsq += temp*temp/(bins1[j]+bins2[j]);
        }
    *prob=gammq(0.5*(*df),0.5*(*chsq));                 Chi-square probability function. See §6.2.
}
```

Equation (14.3.2) and the routine `chstwo` both apply to the case where the total number of data points is the same in the two binned sets. For unequal numbers of data points, the formula analogous to (14.3.2) is

$$\chi^2 = \sum_i \frac{(\sqrt{S/R}R_i - \sqrt{R/S}S_i)^2}{R_i + S_i} \quad (14.3.3)$$

where

$$R \equiv \sum_i R_i \quad S \equiv \sum_i S_i \quad (14.3.4)$$

are the respective numbers of data points. It is straightforward to make the corresponding change in `chstwo`.

Kolmogorov-Smirnov Test

The Kolmogorov-Smirnov (or *K-S*) test is applicable to unbinned distributions that are functions of a single independent variable, that is, to data sets where each data point can be associated with a single number (lifetime of each lightbulb when it burns out, or declination of each star). In such cases, the list of data points can be easily converted to an unbiased estimator $S_N(x)$ of the *cumulative* distribution function of the probability distribution from which it was drawn: If the N events are located at values x_i , $i = 1, \dots, N$, then $S_N(x)$ is the function giving the fraction of data points to the left of a given value x . This function is obviously constant between consecutive (i.e., sorted into ascending order) x_i 's, and jumps by the same constant $1/N$ at each x_i . (See Figure 14.3.1.)

Different distribution functions, or sets of data, give different cumulative distribution function estimates by the above procedure. However, all cumulative distribution functions agree at the smallest allowable value of x (where they are zero), and at the largest allowable value of x (where they are unity). (The smallest and largest values might of course be $\pm\infty$.) So it is the behavior between the largest and smallest values that distinguishes distributions.

One can think of any number of statistics to measure the overall difference between two cumulative distribution functions: the absolute value of the area between them, for example. Or their integrated mean square difference. The Kolmogorov-Smirnov D is a particularly simple measure: It is defined as the *maximum value* of the absolute difference between two cumulative distribution functions. Thus, for comparing one data set's $S_N(x)$ to a known cumulative distribution function $P(x)$, the K-S statistic is

$$D = \max_{-\infty < x < \infty} |S_N(x) - P(x)| \quad (14.3.5)$$

while for comparing two different cumulative distribution functions $S_{N_1}(x)$ and $S_{N_2}(x)$, the K-S statistic is

$$D = \max_{-\infty < x < \infty} |S_{N_1}(x) - S_{N_2}(x)| \quad (14.3.6)$$

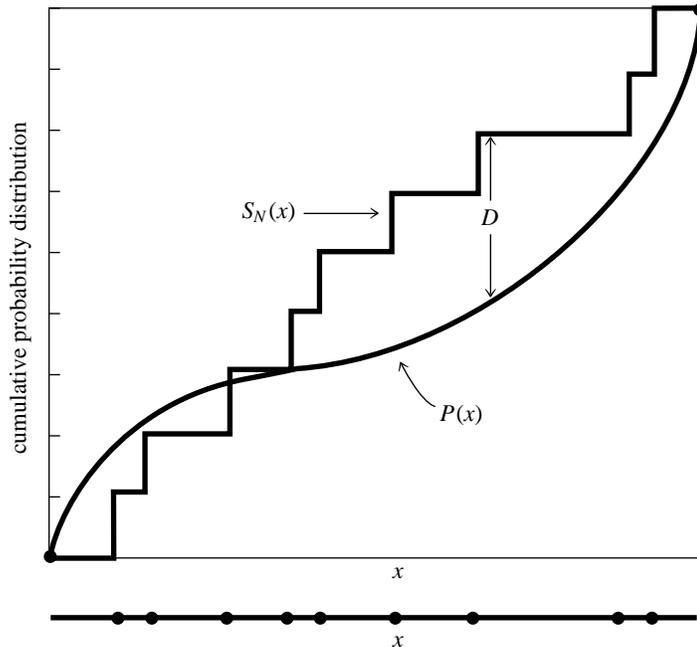


Figure 14.3.1. Kolmogorov-Smirnov statistic D . A measured distribution of values in x (shown as N dots on the lower abscissa) is to be compared with a theoretical distribution whose cumulative probability distribution is plotted as $P(x)$. A step-function cumulative probability distribution $S_N(x)$ is constructed, one that rises an equal amount at each measured point. D is the greatest distance between the two cumulative distributions.

What makes the K–S statistic useful is that *its* distribution in the case of the null hypothesis (data sets drawn from the same distribution) can be calculated, at least to useful approximation, thus giving the significance of any observed nonzero value of D . A central feature of the K–S test is that it is invariant under reparametrization of x ; in other words, you can locally slide or stretch the x axis in Figure 14.3.1, and the maximum distance D remains unchanged. For example, you will get the same significance using x as using $\log x$.

The function that enters into the calculation of the significance can be written as the following sum:

$$Q_{KS}(\lambda) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} e^{-2j^2\lambda^2} \quad (14.3.7)$$

which is a monotonic function with the limiting values

$$Q_{KS}(0) = 1 \quad Q_{KS}(\infty) = 0 \quad (14.3.8)$$

In terms of this function, the significance level of an observed value of D (as a disproof of the null hypothesis that the distributions are the same) is given approximately [1] by the formula

$$\text{Probability}(D > \text{observed}) = Q_{KS}\left(\left[\sqrt{N_e} + 0.12 + 0.11/\sqrt{N_e}\right] D\right) \quad (14.3.9)$$

where N_e is the effective number of data points, $N_e = N$ for the case (14.3.5) of one distribution, and

$$N_e = \frac{N_1 N_2}{N_1 + N_2} \quad (14.3.10)$$

for the case (14.3.6) of two distributions, where N_1 is the number of data points in the first distribution, N_2 the number in the second.

The nature of the approximation involved in (14.3.9) is that it becomes asymptotically accurate as the N_e becomes large, but is already quite good for $N_e \geq 4$, as small a number as one might ever actually use. (See [1].)

So, we have the following routines for the cases of one and two distributions:

```
#include <math.h>
#include "nrutil.h"

void ksone(float data[], unsigned long n, float (*func)(float), float *d,
           float *prob)
Given an array data[1..n], and given a user-supplied function of a single variable func which
is a cumulative distribution function ranging from 0 (for smallest values of its argument) to 1
(for largest values of its argument), this routine returns the K-S statistic d, and the significance
level prob. Small values of prob show that the cumulative distribution function of data is
significantly different from func. The array data is modified by being sorted into ascending
order.
{
    float probks(float alam);
    void sort(unsigned long n, float arr[]);
    unsigned long j;
    float dt,en,ff,fn,fo=0.0;

    sort(n,data);
    en=n;
    *d=0.0;
    for (j=1;j<=n;j++) {
        fn=j/en;
        ff=(*func)(data[j]);
        dt=FMAX(fabs(fo-ff),fabs(fn-ff));
        if (dt > *d) *d=dt;
        fo=fn;
    }
    en=sqrt(en);
    *prob=probks((en+0.12+0.11/en)*(*d));
}

If the data are already sorted into ascending order, then this call can be omitted.
Loop over the sorted data points.
Data's c.d.f. after this step.
Compare to the user-supplied function.
Maximum distance.
Compute significance.

#include <math.h>

void kstwo(float data1[], unsigned long n1, float data2[], unsigned long n2,
           float *d, float *prob)
Given an array data1[1..n1], and an array data2[1..n2], this routine returns the K-S
statistic d, and the significance level prob for the null hypothesis that the data sets are
drawn from the same distribution. Small values of prob show that the cumulative distribution
function of data1 is significantly different from that of data2. The arrays data1 and data2
are modified by being sorted into ascending order.
{
    float probks(float alam);
    void sort(unsigned long n, float arr[]);
    unsigned long j1=1,j2=1;
    float d1,d2,dt,en1,en2,en,fn1=0.0,fn2=0.0;
```

```

sort(n1,data1);
sort(n2,data2);
en1=n1;
en2=n2;
*d=0.0;
while (j1 <= n1 && j2 <= n2) {
    if ((d1=data1[j1]) <= (d2=data2[j2])) fn1=j1++/en1;
    if (d2 <= d1) fn2=j2++/en2;
    if ((dt=fabs(fn2-fn1)) > *d) *d=dt;
}
en=sqrt(en1*en2/(en1+en2));
*prob=probks((en+0.12+0.11/en)*(*d));
}

```

If we are not done...
Next step is in data1.
Next step is in data2.

Compute significance.

Both of the above routines use the following routine for calculating the function

Q_{KS} :

```

#include <math.h>
#define EPS1 0.001
#define EPS2 1.0e-8

float probks(float alam)
Kolmogorov-Smirnov probability function.
{
    int j;
    float a2,fac=2.0,sum=0.0,term,termbf=0.0;

    a2 = -2.0*alam*alam;
    for (j=1;j<=100;j++) {
        term=fac*exp(a2*j*j);
        sum += term;
        if (fabs(term) <= EPS1*termbf || fabs(term) <= EPS2*sum) return sum;
        fac = -fac;           Alternating signs in sum.
        termbf=fabs(term);
    }
    return 1.0;           Get here only by failing to converge.
}

```

Variants on the K-S Test

The sensitivity of the K-S test to deviations from a cumulative distribution function $P(x)$ is not independent of x . In fact, the K-S test tends to be most sensitive around the median value, where $P(x) = 0.5$, and less sensitive at the extreme ends of the distribution, where $P(x)$ is near 0 or 1. The reason is that the difference $|S_N(x) - P(x)|$ does not, in the null hypothesis, have a probability distribution that is independent of x . Rather, its variance is proportional to $P(x)[1 - P(x)]$, which is largest at $P = 0.5$. Since the K-S statistic (14.3.5) is the maximum difference over all x of two cumulative distribution functions, a deviation that might be statistically significant at *its own* value of x gets compared to the expected chance deviation at $P = 0.5$, and is thus discounted. A result is that, while the K-S test is good at finding *shifts* in a probability distribution, especially changes in the median value, it is not always so good at finding *spreads*, which more affect the tails of the probability distribution, and which may leave the median unchanged.

One way of increasing the power of the K-S statistic out on the tails is to replace D (equation 14.3.5) by a so-called *stabilized* or *weighted* statistic [2-4], for example the *Anderson-Darling statistic*,

$$D^* = \max_{-\infty < x < \infty} \frac{|S_N(x) - P(x)|}{\sqrt{P(x)[1 - P(x)]}} \quad (14.3.11)$$

Unfortunately, there is no simple formula analogous to equations (14.3.7) and (14.3.9) for this statistic, although Noé [5] gives a computational method using a recursion relation and provides a graph of numerical results. There are many other possible similar statistics, for example

$$D^{**} = \int_{P=0}^1 \frac{|S_N(x) - P(x)|}{\sqrt{P(x)[1 - P(x)]}} dP(x) \quad (14.3.12)$$

which is also discussed by Anderson and Darling (see [3]).

Another approach, which we prefer as simpler and more direct, is due to Kuiper [6,7]. We already mentioned that the standard K-S test is invariant under reparametrizations of the variable x . An even more general symmetry, which guarantees equal sensitivities at all values of x , is to wrap the x axis around into a circle (identifying the points at $\pm\infty$), and to look for a statistic that is now invariant under all shifts and parametrizations on the circle. This allows, for example, a probability distribution to be “cut” at some central value of x , and the left and right halves to be interchanged, without altering the statistic or its significance.

Kuiper’s statistic, defined as

$$V = D_+ + D_- = \max_{-\infty < x < \infty} [S_N(x) - P(x)] + \max_{-\infty < x < \infty} [P(x) - S_N(x)] \quad (14.3.13)$$

is the sum of the maximum distance of $S_N(x)$ above and below $P(x)$. You should be able to convince yourself that this statistic has the desired invariance on the circle: Sketch the indefinite integral of two probability distributions defined on the circle as a function of angle around the circle, as the angle goes through several times 360° . If you change the starting point of the integration, D_+ and D_- change individually, but their sum is constant.

Furthermore, there is a simple formula for the asymptotic distribution of the statistic V , directly analogous to equations (14.3.7)–(14.3.10). Let

$$Q_{KP}(\lambda) = 2 \sum_{j=1}^{\infty} (4j^2 \lambda^2 - 1) e^{-2j^2 \lambda^2} \quad (14.3.14)$$

which is monotonic and satisfies

$$Q_{KP}(0) = 1 \quad Q_{KP}(\infty) = 0 \quad (14.3.15)$$

In terms of this function the significance level is [1]

$$\text{Probability}(V > \text{observed}) = Q_{KP} \left(\left[\sqrt{N_e} + 0.155 + 0.24/\sqrt{N_e} \right] D \right) \quad (14.3.16)$$

Here N_e is N in the one-sample case, or is given by equation (14.3.10) in the case of two samples.

Of course, Kuiper’s test is ideal for any problem originally defined on a circle, for example, to test whether the distribution in longitude of something agrees with some theory, or whether two somethings have different distributions in longitude. (See also [8].)

We will leave to you the coding of routines analogous to `ksone`, `kstwo`, and `probks`, above. (For $\lambda < 0.4$, don’t try to do the sum 14.3.14. Its value is 1, to 7 figures, but the series can require many terms to converge, and loses accuracy to roundoff.)

Two final cautionary notes: First, we should mention that all varieties of K-S test lack the ability to discriminate some kinds of distributions. A simple example is a probability distribution with a narrow “notch” within which the probability falls to zero. Such a distribution is of course ruled out by the existence of even one data point within the notch, but, because of its cumulative nature, a K-S test would require many data points in the notch before signaling a discrepancy.

Second, we should note that, if you estimate any parameters from a data set (e.g., a mean and variance), then the distribution of the K-S statistic D for a cumulative distribution function $P(x)$ that uses the estimated parameters is no longer given by equation (14.3.9). In general, you will have to determine the new distribution yourself, e.g., by Monte Carlo methods.

CITED REFERENCES AND FURTHER READING:

von Mises, R. 1964, *Mathematical Theory of Probability and Statistics* (New York: Academic Press), Chapters IX(C) and IX(E).

- Stephens, M.A. 1970, *Journal of the Royal Statistical Society*, ser. B, vol. 32, pp. 115–122. [1]
- Anderson, T.W., and Darling, D.A. 1952, *Annals of Mathematical Statistics*, vol. 23, pp. 193–212. [2]
- Darling, D.A. 1957, *Annals of Mathematical Statistics*, vol. 28, pp. 823–838. [3]
- Michael, J.R. 1983, *Biometrika*, vol. 70, no. 1, pp. 11–17. [4]
- Noé, M. 1972, *Annals of Mathematical Statistics*, vol. 43, pp. 58–64. [5]
- Kuiper, N.H. 1962, *Proceedings of the Koninklijke Nederlandse Akademie van Wetenschappen*, ser. A., vol. 63, pp. 38–47. [6]
- Stephens, M.A. 1965, *Biometrika*, vol. 52, pp. 309–321. [7]
- Fisher, N.I., Lewis, T., and Embleton, B.J.J. 1987, *Statistical Analysis of Spherical Data* (New York: Cambridge University Press). [8]

14.4 Contingency Table Analysis of Two Distributions

In this section, and the next two sections, we deal with *measures of association* for two distributions. The situation is this: Each data point has two or more different quantities associated with it, and we want to know whether knowledge of one quantity gives us any demonstrable advantage in predicting the value of another quantity. In many cases, one variable will be an “independent” or “control” variable, and another will be a “dependent” or “measured” variable. Then, we want to know if the latter variable *is* in fact dependent on or *associated* with the former variable. If it is, we want to have some quantitative measure of the strength of the association. One often hears this loosely stated as the question of whether two variables are *correlated* or *uncorrelated*, but we will reserve those terms for a particular kind of association (linear, or at least monotonic), as discussed in §14.5 and §14.6.

Notice that, as in previous sections, the different concepts of significance and strength appear: The association between two distributions may be very significant even if that association is weak — if the quantity of data is large enough.

It is useful to distinguish among some different kinds of variables, with different categories forming a loose hierarchy.

- A variable is called *nominal* if its values are the members of some unordered set. For example, “state of residence” is a nominal variable that (in the U.S.) takes on one of 50 values; in astrophysics, “type of galaxy” is a nominal variable with the three values “spiral,” “elliptical,” and “irregular.”
- A variable is termed *ordinal* if its values are the members of a discrete, but ordered, set. Examples are: grade in school, planetary order from the Sun (Mercury = 1, Venus = 2, . . .), number of offspring. There need not be any concept of “equal metric distance” between the values of an ordinal variable, only that they be intrinsically ordered.
- We will call a variable *continuous* if its values are real numbers, as are times, distances, temperatures, etc. (Social scientists sometimes distinguish between *interval* and *ratio* continuous variables, but we do not find that distinction very compelling.)