

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.)

A continuous variable can always be made into an ordinal one by binning it into ranges. If we choose to ignore the ordering of the bins, then we can turn it into a nominal variable. Nominal variables constitute the lowest type of the hierarchy, and therefore the most general. For example, a set of *several* continuous or ordinal variables can be turned, if crudely, into a single nominal variable, by coarsely binning each variable and then taking each distinct combination of bin assignments as a single nominal value. When multidimensional data are sparse, this is often the only sensible way to proceed.

The remainder of this section will deal with measures of association between *nominal* variables. For any pair of nominal variables, the data can be displayed as a *contingency table*, a table whose rows are labeled by the values of one nominal variable, whose columns are labeled by the values of the other nominal variable, and whose entries are nonnegative integers giving the number of observed events for each combination of row and column (see Figure 14.4.1). The analysis of association between nominal variables is thus called *contingency table analysis* or *cross-tabulation analysis*.

We will introduce two different approaches. The first approach, based on the chi-square statistic, does a good job of characterizing the significance of association, but is only so-so as a measure of the strength (principally because its numerical values have no very direct interpretations). The second approach, based on the information-theoretic concept of *entropy*, says nothing at all about the significance of association (use chi-square for that!), but is capable of very elegantly characterizing the strength of an association already known to be significant.

Measures of Association Based on Chi-Square

Some notation first: Let N_{ij} denote the number of events that occur with the first variable x taking on its i th value, and the second variable y taking on its j th value. Let N denote the total number of events, the sum of all the N_{ij} 's. Let $N_{i\cdot}$ denote the number of events for which the first variable x takes on its i th value regardless of the value of y ; $N_{\cdot j}$ is the number of events with the j th value of y regardless of x . So we have

$$\begin{aligned} N_{i\cdot} &= \sum_j N_{ij} & N_{\cdot j} &= \sum_i N_{ij} \\ N &= \sum_i N_{i\cdot} = \sum_j N_{\cdot j} \end{aligned} \tag{14.4.1}$$

	1. red	2. green	...	
1. male	# of red males N_{11}	# of green males N_{12}	...	# of males $N_{1\cdot}$
2. female	# of red females N_{21}	# of green females N_{22}	...	# of females $N_{2\cdot}$
⋮	⋮	⋮	⋮	⋮
	# of red $N_{\cdot 1}$	# of green $N_{\cdot 2}$...	total # N

Figure 14.4.1. Example of a contingency table for two nominal variables, here sex and color. The row and column marginals (totals) are shown. The variables are “nominal,” i.e., the order in which their values are listed is arbitrary and does not affect the result of the contingency table analysis. If the ordering of values has some intrinsic meaning, then the variables are “ordinal” or “continuous,” and correlation techniques (§14.5-§14.6) can be utilized.

$N_{\cdot j}$ and $N_{i\cdot}$ are sometimes called the *row and column totals* or *marginals*, but we will use these terms cautiously since we can never keep straight which are the rows and which are the columns!

The null hypothesis is that the two variables x and y have no association. In this case, the probability of a particular value of x given a particular value of y should be the same as the probability of that value of x regardless of y . Therefore, in the null hypothesis, the expected number for any N_{ij} , which we will denote n_{ij} , can be calculated from only the row and column totals,

$$\frac{n_{ij}}{N_{\cdot j}} = \frac{N_{i\cdot}}{N} \quad \text{which implies} \quad n_{ij} = \frac{N_{i\cdot} \cdot N_{\cdot j}}{N} \quad (14.4.2)$$

Notice that if a column or row total is zero, then the expected number for all the entries in that column or row is also zero; in that case, the never-occurring bin of x or y should simply be removed from the analysis.

The chi-square statistic is now given by equation (14.3.1), which, in the present case, is summed over all entries in the table,

$$\chi^2 = \sum_{i,j} \frac{(N_{ij} - n_{ij})^2}{n_{ij}} \quad (14.4.3)$$

The number of degrees of freedom is equal to the number of entries in the table (product of its row size and column size) minus the number of constraints that have arisen from our use of the data themselves to determine the n_{ij} . Each row total and column total is a constraint, except that this overcounts by one, since the total of the

column totals and the total of the row totals both equal N , the total number of data points. Therefore, if the table is of size I by J , the number of degrees of freedom is $IJ - I - J + 1$. Equation (14.4.3), along with the chi-square probability function (§6.2), now give the significance of an association between the variables x and y .

Suppose there is a significant association. How do we quantify its strength, so that (e.g.) we can compare the strength of one association with another? The idea here is to find some reparametrization of χ^2 which maps it into some convenient interval, like 0 to 1, where the result is not dependent on the quantity of data that we happen to sample, but rather depends only on the underlying population from which the data were drawn. There are several different ways of doing this. Two of the more common are called *Cramer's V* and the *contingency coefficient C*.

The formula for Cramer's V is

$$V = \sqrt{\frac{\chi^2}{N \min(I-1, J-1)}} \quad (14.4.4)$$

where I and J are again the numbers of rows and columns, and N is the total number of events. Cramer's V has the pleasant property that it lies between zero and one inclusive, equals zero when there is no association, and equals one only when the association is perfect: All the events in any row lie in one unique column, and vice versa. (In chess parlance, no two rooks, placed on a nonzero table entry, can capture each other.)

In the case of $I = J = 2$, Cramer's V is also referred to as the *phi* statistic.

The contingency coefficient C is defined as

$$C = \sqrt{\frac{\chi^2}{\chi^2 + N}} \quad (14.4.5)$$

It also lies between zero and one, but (as is apparent from the formula) it can never achieve the upper limit. While it can be used to compare the strength of association of two tables with the same I and J , its upper limit depends on I and J . Therefore it can never be used to compare tables of different sizes.

The trouble with both Cramer's V and the contingency coefficient C is that, when they take on values in between their extremes, there is no very direct interpretation of what that value means. For example, you are in Las Vegas, and a friend tells you that there is a small, but significant, association between the color of a croupier's eyes and the occurrence of red and black on his roulette wheel. Cramer's V is about 0.028, your friend tells you. You know what the usual odds against you are (because of the green zero and double zero on the wheel). Is this association sufficient for you to make money? Don't ask us!

```

SUBROUTINE cntab1(nn,ni,nj, chisq,df,prob,cramrv,ccc)
INTEGER ni,nj,nn(ni,nj),MAXI,MAXJ
REAL ccc,chisq,cramrv,df,prob,TINY
PARAMETER (MAXI=100,MAXJ=100,TINY=1.e-30)      Maximum table size, and a small number.
C  USES gammq
    Given a two-dimensional contingency table in the form of an integer array nn(1:ni,1:nj),
    this routine returns the chi-square chisq, the number of degrees of freedom df, the significance
    level prob (small values indicating a significant association), and two measures of
    association, Cramer's V (cramrv) and the contingency coefficient C (ccc).
INTEGER i,j,nni,nnj

```

```

REAL expctd,sum,sumi(MAXI),sumj(MAXJ),gammq
sum=0
nni=ni
nnj=nj
do 12 i=1,ni
  sumi(i)=0.
  do 11 j=1,nj
    sumi(i)=sumi(i)+nn(i,j)
    sum=sum+nn(i,j)
  enddo 11
  if(sumi(i).eq.0.)nni=nni-1
enddo 12
do 14 j=1,nj
  sumj(j)=0.
  do 13 i=1,ni
    sumj(j)=sumj(j)+nn(i,j)
  enddo 13
  if(sumj(j).eq.0.)nnj=nnj-1
enddo 14
df=nni*nnj-nni-nnj+1
chisq=0.
do 16 i=1,ni
  do 15 j=1,nj
    expctd=sumj(j)*sumi(i)/sum
    chisq=chisq+(nn(i,j)-expctd)**2/(expctd+TINY)
  enddo 15
enddo 16
prob=gammq(0.5*df,0.5*chisq)
cramrv=sqrt(chisq/(sum*min(nni-1,nnj-1)))
ccc=sqrt(chisq/(chisq+sum))
return
END

```

Will be total number of events.
Number of rows
and columns.
Get the row totals.
Eliminate any zero rows by reducing the
number.
Get the column totals.
Eliminate any zero columns.
Corrected number of degrees of freedom.
Do the chi-square sum.
Here TINY guarantees that
any eliminated row or column will not
contribute to the sum.
Chi-square probability function.

Measures of Association Based on Entropy

Consider the game of “twenty questions,” where by repeated yes/no questions you try to eliminate all except one correct possibility for an unknown object. Better yet, consider a generalization of the game, where you are allowed to ask multiple choice questions as well as binary (yes/no) ones. The categories in your multiple choice questions are supposed to be mutually exclusive and exhaustive (as are “yes” and “no”).

The value to you of an answer increases with the number of possibilities that it eliminates. More specifically, an answer that eliminates all except a fraction p of the remaining possibilities can be assigned a value $-\ln p$ (a positive number, since $p < 1$). The purpose of the logarithm is to make the value additive, since (e.g.) one question that eliminates all but 1/6 of the possibilities is considered as good as two questions that, in sequence, reduce the number by factors 1/2 and 1/3.

So that is the value of an answer; but what is the value of a question? If there are I possible answers to the question ($i = 1, \dots, I$) and the fraction of possibilities consistent with the i th answer is p_i (with the sum of the p_i 's equal to one), then the value of the question is the expectation value of the value of the answer, denoted H ,

$$H = - \sum_{i=1}^I p_i \ln p_i \quad (14.4.6)$$

In evaluating (14.4.6), note that

$$\lim_{p \rightarrow 0} p \ln p = 0 \quad (14.4.7)$$

The value H lies between 0 and $\ln I$. It is zero only when one of the p_i 's is one, all the others zero: In this case, the question is valueless, since its answer is preordained. H takes on its maximum value when all the p_i 's are equal, in which case the question is sure to eliminate all but a fraction $1/I$ of the remaining possibilities.

The value H is conventionally termed the *entropy* of the distribution given by the p_i 's, a terminology borrowed from statistical physics.

So far we have said nothing about the association of two variables; but suppose we are deciding what question to ask next in the game and have to choose between two candidates, or possibly want to ask both in one order or another. Suppose that one question, x , has I possible answers, labeled by i , and that the other question, y , as J possible answers, labeled by j . Then the possible outcomes of asking both questions form a contingency table whose entries N_{ij} , when normalized by dividing by the total number of remaining possibilities N , give all the information about the p 's. In particular, we can make contact with the notation (14.4.1) by identifying

$$\begin{aligned} p_{ij} &= \frac{N_{ij}}{N} \\ p_{i.} &= \frac{N_{i.}}{N} \quad (\text{outcomes of question } x \text{ alone}) \\ p_{.j} &= \frac{N_{.j}}{N} \quad (\text{outcomes of question } y \text{ alone}) \end{aligned} \quad (14.4.8)$$

The entropies of the questions x and y are, respectively,

$$H(x) = - \sum_i p_{i.} \ln p_{i.} \quad H(y) = - \sum_j p_{.j} \ln p_{.j} \quad (14.4.9)$$

The entropy of the two questions together is

$$H(x, y) = - \sum_{i,j} p_{ij} \ln p_{ij} \quad (14.4.10)$$

Now what is the entropy of the question y given x (that is, if x is asked first)? It is the expectation value over the answers to x of the entropy of the restricted y distribution that lies in a single column of the contingency table (corresponding to the x answer):

$$H(y|x) = - \sum_i p_{i.} \sum_j \frac{p_{ij}}{p_{i.}} \ln \frac{p_{ij}}{p_{i.}} = - \sum_{i,j} p_{ij} \ln \frac{p_{ij}}{p_{i.}} \quad (14.4.11)$$

Correspondingly, the entropy of x given y is

$$H(x|y) = - \sum_j p_{.j} \sum_i \frac{p_{ij}}{p_{.j}} \ln \frac{p_{ij}}{p_{.j}} = - \sum_{i,j} p_{ij} \ln \frac{p_{ij}}{p_{.j}} \quad (14.4.12)$$

We can readily prove that the entropy of y given x is never more than the entropy of y alone, i.e., that asking x first can only reduce the usefulness of asking y (in which case the two variables are *associated!*):

$$\begin{aligned}
 H(y|x) - H(y) &= - \sum_{i,j} p_{ij} \ln \frac{p_{ij}/p_{i\cdot}}{p_{\cdot j}} \\
 &= \sum_{i,j} p_{ij} \ln \frac{p_{\cdot j} p_{i\cdot}}{p_{ij}} \\
 &\leq \sum_{i,j} p_{ij} \left(\frac{p_{\cdot j} p_{i\cdot}}{p_{ij}} - 1 \right) \\
 &= \sum_{i,j} p_{i\cdot} p_{\cdot j} - \sum_{i,j} p_{ij} \\
 &= 1 - 1 = 0
 \end{aligned} \tag{14.4.13}$$

where the inequality follows from the fact

$$\ln w \leq w - 1 \tag{14.4.14}$$

We now have everything we need to define a measure of the “dependency” of y on x , that is to say a measure of association. This measure is sometimes called the *uncertainty coefficient* of y . We will denote it as $U(y|x)$,

$$U(y|x) \equiv \frac{H(y) - H(y|x)}{H(y)} \tag{14.4.15}$$

This measure lies between zero and one, with the value 0 indicating that x and y have no association, the value 1 indicating that knowledge of x completely predicts y . For in-between values, $U(y|x)$ gives the fraction of y 's entropy $H(y)$ that is lost if x is already known (i.e., that is redundant with the information in x). In our game of “twenty questions,” $U(y|x)$ is the fractional loss in the utility of question y if question x is to be asked first.

If we wish to view x as the dependent variable, y as the independent one, then interchanging x and y we can of course define the dependency of x on y ,

$$U(x|y) \equiv \frac{H(x) - H(x|y)}{H(x)} \tag{14.4.16}$$

If we want to treat x and y symmetrically, then the useful combination turns out to be

$$U(x, y) \equiv 2 \left[\frac{H(y) + H(x) - H(x, y)}{H(x) + H(y)} \right] \tag{14.4.17}$$

If the two variables are completely independent, then $H(x, y) = H(x) + H(y)$, so (14.4.17) vanishes. If the two variables are completely dependent, then $H(x) = H(y) = H(x, y)$, so (14.4.16) equals unity. In fact, you can use the identities (easily proved from equations 14.4.9–14.4.12)

$$H(x, y) = H(x) + H(y|x) = H(y) + H(x|y) \tag{14.4.18}$$

to show that

$$U(x, y) = \frac{H(x)U(x|y) + H(y)U(y|x)}{H(x) + H(y)} \quad (14.4.19)$$

i.e., that the symmetrical measure is just a weighted average of the two asymmetrical measures (14.4.15) and (14.4.16), weighted by the entropy of each variable separately.

Here is a program for computing all the quantities discussed, $H(x)$, $H(y)$, $H(x|y)$, $H(y|x)$, $H(x, y)$, $U(x|y)$, $U(y|x)$, and $U(x, y)$:

```

SUBROUTINE cntab2(nn,ni,nj,h,hx,hy,hygx,hxgy,uygx,uxgy,uxy)
INTEGER ni,nj,nn(ni,nj),MAXI,MAXJ
REAL h,hx,hxgy,hy,hygx,uxgy,uxy,uygx,TINY
PARAMETER (MAXI=100,MAXJ=100,TINY=1.e-30)
  Given a two-dimensional contingency table in the form of an integer array nn(i,j), where
  i labels the x variable and ranges from 1 to ni, j labels the y variable and ranges from 1 to
  nj, this routine returns the entropy h of the whole table, the entropy hx of the x distribution,
  the entropy hy of the y distribution, the entropy hygx of y given x, the entropy hxgy of
  x given y, the dependency uygx of y on x (eq. 14.4.15), the dependency uxgy of x on y
  (eq. 14.4.16), and the symmetrical dependency uxy (eq. 14.4.17).
  Parameters: MAXI and MAXJ define the maximum size of table; TINY is a small number.
  INTEGER i,j
  REAL p,sum,sumi(MAXI),sumj(MAXJ)
  sum=0
  do 12 i=1,ni
    sumi(i)=0.0
    do 11 j=1,nj
      sumi(i)=sumi(i)+nn(i,j)
      sum=sum+nn(i,j)
    enddo 11
  enddo 12
  do 14 j=1,nj
    sumj(j)=0.
    do 13 i=1,ni
      sumj(j)=sumj(j)+nn(i,j)
    enddo 13
  enddo 14
  hx=0.
  do 15 i=1,ni
    if(sumi(i).ne.0.)then
      p=sumi(i)/sum
      hx=hx-p*log(p)
    endif
  enddo 15
  hy=0.
  do 16 j=1,nj
    if(sumj(j).ne.0.)then
      p=sumj(j)/sum
      hy=hy-p*log(p)
    endif
  enddo 16
  h=0.
  do 18 i=1,ni
    do 17 j=1,nj
      if(nn(i,j).ne.0)then
        p=nn(i,j)/sum
        h=h-p*log(p)
      endif
    enddo 17
  enddo 18
  hygx=h-hx
  hxgy=h-hy

```

Get the row totals.

Get the column totals.

Entropy of the x distribution,

and of the y distribution.

Total entropy: loop over both x and y.

Uses equation (14.4.18), as does this.


```

uygx=(hy-hygx)/(hy+TINY)           Equation (14.4.15).
uxgy=(hx-hxgy)/(hx+TINY)           Equation (14.4.16).
uxy=2.*(hx+hy-h)/(hx+hy+TINY)      Equation (14.4.17).
return
END

```

CITED REFERENCES AND FURTHER READING:

- Dunn, O.J., and Clark, V.A. 1974, *Applied Statistics: Analysis of Variance and Regression* (New York: Wiley).
- Norusis, M.J. 1982, *SPSS Introductory Guide: Basic Statistics and Operations*, and 1985, *SPSS-X Advanced Statistics Guide* (New York: McGraw-Hill).
- Fano, R.M. 1961, *Transmission of Information* (New York: Wiley and MIT Press), Chapter 2.

14.5 Linear Correlation

We next turn to measures of association between variables that are ordinal or continuous, rather than nominal. Most widely used is the *linear correlation coefficient*. For pairs of quantities (x_i, y_i) , $i = 1, \dots, N$, the linear correlation coefficient r (also called the product-moment correlation coefficient, or *Pearson's r*) is given by the formula

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}} \quad (14.5.1)$$

where, as usual, \bar{x} is the mean of the x_i 's, \bar{y} is the mean of the y_i 's.

The value of r lies between -1 and 1 , inclusive. It takes on a value of 1 , termed "complete positive correlation," when the data points lie on a perfect straight line with positive slope, with x and y increasing together. The value 1 holds independent of the magnitude of the slope. If the data points lie on a perfect straight line with negative slope, y decreasing as x increases, then r has the value -1 ; this is called "complete negative correlation." A value of r near zero indicates that the variables x and y are *uncorrelated*.

When a correlation is known to be significant, r is one conventional way of summarizing its strength. In fact, the value of r can be translated into a statement about what residuals (root mean square deviations) are to be expected if the data are fitted to a straight line by the least-squares method (see §15.2, especially equations 15.2.13 – 15.2.14). Unfortunately, r is a rather poor statistic for deciding *whether* an observed correlation is statistically significant, and/or whether one observed correlation is significantly stronger than another. The reason is that r is ignorant of the individual distributions of x and y , so there is no universal way to compute its distribution in the case of the null hypothesis.

About the only general statement that can be made is this: If the null hypothesis is that x and y are uncorrelated, and if the distributions for x and y each have enough convergent moments ("tails" die off sufficiently rapidly), and if N is large