

Categorical Data Analysis

→ Categorical Variable and Categorical Data: ~ A variable which takes values on a scale consisting of some categories is called a categorical variable. For example, "Political Ideology" is a categorical variable which takes values on a scale consisting of several categories, say, 'liberal', 'moderate' & 'conservative'. Any data collected on a categorical variable is called Categorical Data.

Response & Explanatory Variables: Most statistical techniques distinguish between response or dependent variables and explanatory or independent variables. For instance, regression models describe how the mean of response variable such as the selling prices of houses changes according as the values of explanatory variables such as square footage and location.

Nominal & Ordinal Variables : Categorical variables have two primary types of scales. Variables having categories without a natural ordering are called nominal variables. e.g. religious affiliation with the categories catholic, protestant, jewish, Muslims and others. For nominal variables, the orders of listing of the categories is irrelevant (not necessary). Many categorical variables do have ordered categories, such variables are called ordinal variables. Examples are social class: upper, middle, lower; & patient conditions: good, fair, serious.

Attribute : ~ When we record the sex of each newborn baby during a month or the language of each book in a library, the data are not numbers initially. We get numbers if, subsequently, we note the number of male babies and that of female babies, or the numbers of books written in English, the number written in Hindi, the numbers written in Bengali and so forth. For this type of data, the characters observed is not expressible in numerical terms. Such a character is, therefore, called a qualitative character or an attribute.

$$(S1A), (S2A) : \text{babies born}$$

$$(S3a), (S4a)$$

Qualitative for small of babies and appear to non numerical data

Dichotomy : ~ A classification of the simple kind considered, in which each class is divided into two sub-classes and no more, has been termed by logicians classification, or to use the more strictly applicable term, division by dichotomy; the classification of most statistics are not dichotomous, for most usually a class is divided into more than two subclasses, but dichotomy is the fundamental case.

Notations : ~ For theoretical purposes, it is necessary to have some simple notations for the classes formed and for the numbers of observation assigned to each. Let us denote the several attributes by A, B, C, \dots , etc. An object or an individual possessing the attribute A will be termed simply by A . The class, all the members of which possess the attribute A will be termed as the class- A . The absence of the attribute A , we shall employ the notations $\alpha, \beta, \gamma, \dots$ according to A, B, C, \dots ; we shall represent the attribute blindness; α represents sight. If B stands for deafness, β stands for hearing. Generally, ' α ' is equivalent to 'not A '. The class α is equivalent to the class none of the members of which possesses the attribute A .

In case of combination of attributes, such as, 'AB' represents the combination of blindness and deafness. If a third attribute be noted, e.g., insanity, denoted by C & the absence of it by γ ; then the class 'ABC' includes those who are at once deaf, blind and insane; ' $AB\gamma$ ' includes those who are deaf and blind but not insane.

Class-frequencies : ~ The no. of observation assigned to any class is termed as the 'class-frequency'. Class-frequencies will be denoted by putting the class-symbols within parenthesis.

Orders of classes and class-frequency : ~

For 2 attributes, the class frequency of different orders are

Orders 0 : N

1st orders : $(A), (B), (\alpha), (\beta)$

[(A) denotes the class frequency of the class A]

2nd orders : $(AB), (A\beta), (\alpha B), (\alpha\beta)$

Any class frequency can always be expressed in terms of frequencies of higher orders classes. e.g.

$$(A) = (AB) + (A\beta)$$

$$(AB) = (ABC) + (AB\gamma)$$

Ultimate class frequencies :- A class of highest orders is known as an ultimate class and its frequency an ultimate class frequency.

Every class frequency can be written as a sum of certain ultimate class frequencies.

For any frequency can be analysed into higher frequencies and the process needs stop only when we have reached the frequency of the highest orders, e.g. with 3 attributes

$$(A) = (AB) + (A\bar{B})$$

$$= (ABC) + (AB\bar{C}) + (A\bar{B}C) + (A\bar{B}\bar{C})$$

The total no. of class-frequencies :-

Orders zero : N

Orders one : (A) (B) (C)
 (\bar{A}) (\bar{B}) (\bar{C})

Orders two : (AB) (AC) (BC)
 $(A\bar{B})$ $(A\bar{C})$ $(B\bar{C})$
 $(\bar{A}B)$ $(\bar{A}C)$ $(B\bar{C})$
 $(\bar{A}\bar{B})$ $(\bar{A}\bar{C})$ $(\bar{B}\bar{C})$

Orders three : (ABC) $(\bar{A}\bar{B}\bar{C})$
 $(AB\bar{C})$ $(A\bar{B}C)$ $(B\bar{A}C)$
 $(A\bar{B}\bar{C})$ $(\bar{A}B\bar{C})$
 $(\bar{A}\bar{B}C)$ $(\bar{A}\bar{B}\bar{C})$

Hence we have 3^3 distinct class-frequencies thus, for n attributes there are 3^n distinct in all. In general, for n attributes there are 3^n distinct class-frequencies, provided we count N as a frequency of orders zero.

Of orders zero, there is a single class N ; of order one, there are $2n$ classes; of order two, there are $\binom{n}{2} 2^2$ classes; hence of order n , there are $\binom{n}{n} 2^n$ classes.

$$\therefore \text{No. of class frequencies} = \sum_{n=0}^{\infty} \binom{n}{n} 2^n = (1+2)^n = 3^n$$

Consistency : Any class frequencies which have been observed within one and the same population may be said to be consistent with one another. They conform with one another, and don't, in any way, conflict.

Symbols:

We define, $A \cdot N = (A)$, for any attribute

similarly, $\alpha \cdot N = (\alpha)$

$$\text{Now, } (A + \alpha)N = (A) + (\alpha) = N$$

$$\Rightarrow (A + \alpha) = 1$$

$$\Rightarrow \alpha = 1 - A$$

Therefore, in any symbolic expression we can replace the operators A by $1 - \alpha$ or α by $1 - A$.

$$\text{Again, } (AB) = A \cdot (B) = B \cdot (A)$$

A little reflection will show that the operative symbols therefore obey the laws of algebra.

$$\begin{aligned} \therefore (\alpha\beta) &= \alpha\beta \cdot N \\ &= (1 - A)(1 - B) \cdot N \\ &= (1 - B - A + AB) \cdot N \\ &= N - (A) - (B) + (AB) \end{aligned}$$

$$\text{Similarly, } (\alpha\beta\gamma) = \alpha\beta\gamma \cdot N$$

$$\begin{aligned} &= (1 - A)(1 - B)(1 - C) \cdot N \\ &= (1 - A - B - C + AB + BC + AC - ABC) \cdot N \\ &= N - (A) - (B) - (C) + (AB) + (BC) + (AC) - (ABC) \end{aligned}$$

Condition for Consistency: — The attributes denoted by capitals may be termed positive attributes, on the other hand, attributes denoted by Greek letters, may be termed as negative attributes.

The necessary and sufficient condition for the consistency of a set of independent class-frequencies is that no ultimate class-frequency be negative.

If two attributes are noted, there are 4 ultimate frequencies — (AB) , $(A\bar{B})$, $(\bar{A}B)$, $(\bar{A}\bar{B})$. Expressing them in terms of positive classes, we find the following conditions —

- i) $(AB) \geq 0$
- ii) $(AB) \geq (A) + (B) - N$ [i.e. $(\alpha\beta) \geq 0$]
- iii) $(AB) \leq (A) \Leftrightarrow (A\bar{B}) \geq 0$
- iv) $(AB) \leq (B) \Leftrightarrow (\bar{A}B) \geq 0$

(5)

For 3 attributes the conditions that the 8 ultimate frequencies are not negative will be found to lead to the following —

- i) $(ABC) \geq 0$
- ii) $(ABC) \geq (AB) + (AC) - (A)$
- iii) $(ABC) \geq (AB) + (AC) - (B)$
- iv) $(ABC) \geq (AB) + (AC) - (C)$
- v) $(ABC) \leq (AB)$
- vi) $(ABC) \leq (AC)$
- vii) $(ABC) \leq (BC)$
- viii) $(ABC) \leq N - (A) - (B) - (C) + (AB) + (AC) + (BC)$

Example: Show that neither α nor β can exceed $1/4$ when the followings are given:

$$\frac{(A)}{N} = \alpha, \frac{(B)}{N} = 2\alpha, \frac{(C)}{N} = 3\alpha \text{ & } \frac{(AB)}{N} = \frac{(BC)}{N} = \frac{(CA)}{N} = \beta.$$

Soln. From condition of consistency,

$$\begin{aligned} (AB) &\leq (A) \\ \Rightarrow \beta &\leq \alpha \quad (*) \\ (BC) &\geq (B) + (C) - N \Rightarrow \frac{(BC)}{N} \geq \frac{(B)}{N} + \frac{(C)}{N} - 1 \\ \Rightarrow \beta &\geq 2\alpha + 3\alpha - 1 = 5\alpha - 1 \quad \Leftrightarrow \\ \Rightarrow \beta &\geq 5\alpha - 1 \quad [\text{Using } (*)] \\ \Rightarrow \alpha &> 5\alpha - 1 \\ \Rightarrow \alpha &\leq 1/4 \\ \Rightarrow \beta &\leq \alpha \leq 1/4 \end{aligned}$$

INDEPENDENCE :

Attribute	$\frac{N}{(B)}$	β	$\frac{N}{(A)} - \frac{N}{(B)} =$	TOTAL
A	$(AB) - (A\bar{B})$	$(\bar{A}B)$	(A)	
α	(αB)	$(\alpha \bar{B})$	(α)	
TOTAL	(B)	(β)	N	(αA)

Two related attributes A & B are said to be independent if

$$\begin{aligned} \frac{(AB)}{N} &= \frac{(A\bar{B})}{N} \\ \Leftrightarrow \frac{(AB)}{N} &= \frac{(AB) + (A\bar{B})}{N} = \frac{(A)}{N} \\ \Leftrightarrow \frac{(AB)}{N} &= \frac{(A)}{N} \cdot \frac{(B)}{N} \quad (***) \end{aligned}$$

If the attributes A and B are independent, the proportion of {AB}'s in the population is equal to the proportion of A's multiplied by the proportion of B's. → fundamental rule.

If there is no sort of relationship of any kind between two attributes A and B, we expect to find (the same proportion of A's amongst the B's as amongst the not B's). We have earlier the criterion of independence for A & B -

$$\frac{(AB)}{(B)} = \frac{(A\bar{B})}{(\bar{B})} \quad \text{--- (*)}$$

If this relation holds good, then the following equations must also hold

$$\frac{(\alpha B)}{(B)} = \frac{(\alpha \bar{B})}{(\bar{B})}$$

$$\frac{(A\bar{B})}{(A)} = \frac{(\alpha \bar{B})}{(\alpha)}$$

$$\frac{(A\bar{B})}{(A)} = \frac{(\alpha B)}{(\alpha)}$$

$$(*) \rightarrow \text{gives } \frac{(A\bar{B})}{(B)} = \frac{(A\bar{B})}{(\bar{B})} \quad x \geq y \Leftrightarrow$$

$$\Leftrightarrow \frac{(B) - (AB)}{(B)} = \frac{(B) - (A\bar{B})}{(\bar{B})}$$

$$\text{i.e. } \frac{(\alpha B)}{(B)} = \frac{(\alpha \bar{B})}{(\bar{B})}, \text{ etc.}$$

NOTE: If 2 attributes A & B are independent, then i) α, β are also independent, ii) A, β are also independent.

$$\Rightarrow \text{For, } (\alpha \beta) = N - (A) - (B) + (AB)$$

$$= N - (A) - (B) + \frac{(A)(B)}{N}$$

$$= N - (A) - \frac{(B)}{N}(N - (A))$$

$$= \frac{(A)(N - (A))(N - (B))}{N} \quad (\text{SA})$$

$$(\alpha) = \frac{(\alpha)(\beta)}{N} \quad (\beta)$$

$$\text{ii) } \frac{(AB)}{(B)} = \frac{(AB) + (A\bar{B})}{(B) + (\bar{B})} = \frac{(A)}{N} \quad (\text{SA})$$

$$\Rightarrow (AB) = \frac{(A)(\beta)}{N} \quad \text{if } A \text{ is substituted before out}$$

Interpretation: The advantages of (**) over (*) is that it gives expressions for the second orders frequency in terms of the frequencies of the 1st orders and the total no. of observations alone.

$$\frac{(\beta)}{N} \cdot \frac{(A)}{N} = \frac{(\text{SA})}{N} \Leftrightarrow$$

(SA) for neither odd, frequency nor B from A substitution at E'

The criteria of independence may be expressed in yet a 3rd form, i.e. in terms of 2nd orders frequencies alone.
If A and B are independent then

$$(AB) \cdot (\alpha\beta) = \frac{(A)(B)}{N} \cdot \frac{(\alpha)(\beta)}{N} = \frac{(A)(B)(\alpha)(\beta)}{N^2}$$

$$\Rightarrow (AB) \cdot (\alpha\beta) = (\alpha B)(A\beta)$$

$$\Leftrightarrow \frac{(AB)}{(\alpha B)} = \frac{(A\beta)}{(\alpha\beta)} \quad \text{--- (1)}$$

$$\text{or, } \frac{(AB)}{(A\beta)} = \frac{(\alpha B)}{(\alpha\beta)} \quad \text{--- (2)}$$

(1) may be read : "The ratio of A's to α 's amongst the B's is equal to the ratio of A's to α 's amongst the β 's.

(2) can be interpreted in a similar way.

Example : If the second-order frequencies have the following values, are A and B independent or not?

$$(AB) = 110, (\alpha B) = 90, (A\beta) = 290, (\alpha\beta) = 510$$

Soln. Clearly, $(AB)(\alpha\beta) = 110 \times 510$
 $(A\beta)(\alpha B) = 290 \times 90$

$$\text{so, } (AB)(\alpha\beta) > (A\beta)(\alpha B)$$

So, A and B are not independent.

ASSOCIATION : A and B are said to be positively associated

(or, simply associated) if $(AB) > \frac{(A) \cdot (B)}{N}$

A and B are said to be negatively associated
 (or, simply disassociated) if $(AB) < \frac{(A) \cdot (B)}{N}$

In common language, one speaks of A and B as being 'associated' if they appear together in a number of cases. But here, by positive association between A and B, we mean that they appear together in a greater no. of cases than is to be expected if they are independent.

Complete Association or Disassociation:

The circumstances in which the association between two attributes is said to be complete, can be explained into two cases —

i) We may say that for complete association, all A's must be B's and all B's must be A's, i.e. $[(AB)=0] \& [(\alpha\beta)=0]$ otherwise, we see that A's and B's occurs in the population in equal numbers.

ii) We may adopt a rather wider meaning and say that all A's are B's or all B's are A's, according to whether the A's or the B's are in the minority.

Similarly, complete disassociation may be taken either as the case when no A's are B's and no α 's are β 's, or more widely as the case when either of these statements is true.

Thus two attributes are completely associated if one of them can't occur without the others, though the others may occur without the one.

The symbols $(AB)_0$ and δ : We define $(AB)_0 = \frac{(A)(B)}{N}$, the value of (AB) under the assumption that the attributes are independent.

We shall use the other symbols

$$(\alpha\beta)_0 = \frac{(\alpha)(\beta)}{N}, (\alpha\beta)_0 = \frac{(\alpha)(B)}{N}, (AB)_0 = \frac{(A)(\beta)}{N}$$

If δ denote the excess of (AB) over $(AB)_0$, then keeping the marginal totals fixed, the table reduces to the form —

Attribute	B	β	TOTAL
A	$(AB)_0 + \delta$	$(AB)_0 - \delta$	(A)
α	$(\alpha\beta)_0 - \delta$	$(\alpha\beta)_0 + \delta$	(α)
TOTAL	(B)	(β)	N

Define, $\delta = (AB) - (AB)_0$ and also, quite generally we have —

$$\delta = (AB) - (AB)_0 = (\alpha\beta) - (\alpha\beta)_0 = (AB)_0 - (AB) = (\alpha\beta)_0 - (\alpha\beta)$$

Note that, $\delta = (AB) - (AB)_0 = (AB) - \frac{(A)(B)}{N}$

$$\begin{aligned}\delta &= \frac{N(AB) - (A)(B)}{N} \\ \Rightarrow \delta &= \frac{1}{N} [N(AB) - \{ (AB) + (AB)_0 \} \{ (AB) + (\alpha\beta) \}] \\ &= \frac{1}{N} [\{ (AB) + (AB)_0 + (\alpha\beta) + (\alpha\beta)_0 \} (AB) - \{ (AB) + (AB)_0 \} \{ (AB) + (\alpha\beta) \}] \\ &= \frac{1}{N} [(AB)(\alpha\beta) - (AB)_0(\alpha\beta)_0] \\ \text{Clearly, } \delta &= 0, \quad \delta > 0 \Leftrightarrow (+ve) \text{ association}, \quad \delta < 0 \Leftrightarrow (-ve) \text{ association} \\ \Leftrightarrow A \text{ and } B &\text{ are independent.}\end{aligned}$$

' δ ' determines uniquely the departure from independence. But it may be of interest to measure the intensity of association. It is called ' δ measure' for association between 2 attributes.

Coefficient of association:

$$\delta = \frac{(AB)(\alpha\beta) - (AB)_0(\alpha\beta)_0}{N}$$

Yule's coefficient of association: As a measure of intensity of association between 2 attributes A and B, Yule gave the following coefficient of association:

$$\begin{aligned}Q &= \frac{(AB)(\alpha\beta) - (AB)_0(\alpha\beta)_0}{(AB)(\alpha\beta) + (AB)_0(\alpha\beta)_0} \\ &= \frac{N\delta}{(AB)(\alpha\beta) + (AB)_0(\alpha\beta)_0}\end{aligned}$$

If A and B are independent, $\delta = 0 \Leftrightarrow Q = 0$.

PROPERTIES:

$$\Rightarrow Q = \frac{N\delta}{(AB)(\alpha\beta) + (\alpha\beta + AB)} \text{ increases as } \delta \text{ increases.}$$

* If all the terms containing A are multiplied by a constant, the value of Q remains unaltered. Similar things hold for α , β and B.

Hence, Q is independent of A's and α 's in the data.

\therefore Relative proportion of A's and α 's are $\frac{k(A)}{k(A)+\alpha}$, $\frac{(\alpha)}{k(A)+\alpha}$

$$\text{Now, } Q^* = \frac{k(AB)(\alpha\beta) - k(AB)_0(\alpha\beta)_0}{k(AB)(\alpha\beta) + k(AB)_0(\alpha\beta)_0} = Q$$

Hence, Q is independent of the relative proportion of A's.

N.P. Q is a symmetric measure (coefficient) and it increases from -1 to 1 as the extent of association increases from perfect -ve to perfect +ve. Shown in property 3).

$$\Rightarrow -1 \leq Q \leq 1$$

Proof: $Q = \frac{(AB)(\alpha\beta) - (\alpha B)(A\beta)}{(AB)(\alpha\beta) + (\alpha B)(A\beta)} = \frac{a-b}{a+b}$,
where $a = (AB)(\alpha\beta) \geq 0$ & $b = (\alpha B)(A\beta) \geq 0$

Note that, $Q = 1 - \frac{2b}{a+b} \leq 1$, —①

and, $Q = -1 + \frac{2a}{a+b} \geq -1$ —②

Combining ① & ② $\Rightarrow -1 \leq Q \leq 1$.

Marginal cases: $\Rightarrow Q=0 \Leftrightarrow a=b$ iff A & B are independent.

$Q=1$ iff $b=0$

$\Leftrightarrow (\alpha B)(A\beta)=0$

$\Leftrightarrow (\alpha B)=0$, or, $(A\beta)=0$

$\Leftrightarrow (AB)=(B)$, or, $(AB)=(A)$

\Leftrightarrow A & B are in complete association.

$Q=-1$ iff A & B are in complete disassociation.

c.v. \Rightarrow Yule's coefficient of colligation.

$$\text{Define, } \gamma = \frac{1 - \sqrt{(AB)(\alpha\beta) / (AB)(\alpha\beta)}}{1 + \sqrt{(AB)(\alpha\beta) / (AB)(\alpha\beta)}} = \frac{\sqrt{(AB)(\alpha\beta)} - \sqrt{(AB)(\alpha\beta)}}{\sqrt{(AB)(\alpha\beta)} + \sqrt{(AB)(\alpha\beta)}}$$

Show that $\Rightarrow Q = \frac{2\gamma}{1+\gamma^2}$

Proof: $\gamma = \frac{\sqrt{(AB)(\alpha\beta)} - \sqrt{(AB)(\alpha\beta)}}{\sqrt{(AB)(\alpha\beta)} + \sqrt{(AB)(\alpha\beta)}}$

$$\Rightarrow \frac{1+\gamma}{1-\gamma} = \frac{\sqrt{(AB)(\alpha\beta)}}{\sqrt{(AB)(\alpha\beta)}} = \frac{\sqrt{(AB)(\alpha\beta)}}{\sqrt{(AB)(\alpha\beta)}} = Q$$

$$\Rightarrow \frac{(1+\gamma)^2}{(1-\gamma)^2} = \frac{(AB)(\alpha\beta)}{(AB)(\alpha\beta)}$$

$$\Rightarrow \frac{(1+\gamma)^2 - (1-\gamma)^2}{(1+\gamma)^2 + (1-\gamma)^2} = \frac{(AB)(\alpha\beta) - (AB)(\alpha\beta)}{(AB)(\alpha\beta) + (AB)(\alpha\beta)}$$

$$\Rightarrow \frac{4\gamma}{2(1+\gamma^2)} = Q$$

$$\Rightarrow Q = \frac{2\gamma}{1+\gamma^2}$$

N.P. γ has the same properties as Q , i.e. $-1 \leq \gamma \leq 1$.
holds, i.e. (+) for complete association & (-) for complete disassociation. If A & B are independent, then γ vanishes.

→ Manifold Classification : → Instead of dividing the population under consideration into two parts by a simple dichotomy, we may also divide it into a number of parts by a simple process. For instance, we can extend the dichotomy of the population of men into "those with blue eyes" and "those not with blue eyes" to a threefold division: "those with blue eyes", "those with brown eyes" and "those with neither blue nor brown eyes"; or into a fourfold division by adding a fresh category, "those with grey eyes";

Generally, our population may be divided first according to s heads, A_1, A_2, \dots, A_s ; each of the classes so obtained into t heads, B_1, B_2, \dots, B_t ; each of these into u heads, C_1, C_2, \dots, C_u ; and so on.

This is called "manifold classification".

The theory of manifold classification involving n attributes is rather complicated, but its fundamental principles are very similar to dichotomy; A straightforward extension of the methods already discussed will give the following results:

- (a) There are $sxt \times u \times \dots$ ultimate classes.
- (b) The total no. of classes, including N and the ultimate classes, is $(s+1)(t+1)(u+1) \dots$
- (c) The data are consistent if, and only if, every ultimate class-frequency is not negative.
- (d) The data are completely specified by $sxt \times u \times \dots$ algebraically independent class-frequencies. Even if all these are not given, it may be possible to set limits to the other class-frequencies.

■ Contingency Table : — Let us consider a classification in which the attribute A is s -fold and B is t -fold. Clearly, here we have 'st' classes of the type $A_i B_j$; $i=1(1)s$, $j=1(1)t$.

Then we can arrange our data in a table as follows —

A	A_1	A_2	\dots	A_m	\dots	A_s	Totals
B	$(A_1 B_1)$	$(A_2 B_1)$	\dots	$(A_m B_1)$	\dots	$(A_s B_1)$	(B_1)
B_1	$(A_1 B_2)$	$(A_2 B_2)$	\dots	$(A_m B_2)$	\dots	$(A_s B_2)$	(B_2)
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
B_n	$(A_1 B_n)$	$(A_2 B_n)$	\dots	$(A_m B_n)$	\dots	$(A_s B_n)$	(B_n)
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
B_t	$(A_1 B_t)$	$(A_2 B_t)$	\dots	$(A_m B_t)$	\dots	$(A_s B_t)$	(B_t)
Totals	(A_1)	(A_2)	\dots	(A_m)	\dots	(A_s)	N

Such a table is called a Contingency Table.

Alternative Notations:

$$n_{ij} = (A_i B_j)$$

$$n_{i0} = (A_i)$$

$$n_{0j} = (B_j)$$

Coefficients of Contingency : ~ The characteristics A & B are said to be completely independent in the population at large, if for all values of m and n, we must have —

$$(A_m B_n) = \frac{(A_m)(B_n)}{N} = (A_m B_n)_0, \text{ i.e., } n_{ij} = \frac{n_{i0} \cdot n_{0j}}{N}$$

If, however, A and B are not completely independent, $(A_m B_n)$ and $(A_m B_n)_0$ will not be identical for all the values of m and n. The deviation from independence in that particular cell will be measured by,

$$\delta_{mn} = (A_m B_n) - (A_m B_n)_0 \\ = (A_m B_n) - \frac{(A_m)(B_n)}{N}$$

$$\text{i.e., } \delta_{ij} = n_{ij} - \frac{n_{i0} \cdot n_{0j}}{N}$$

Properties: →

(1) In General $\underline{\delta_{mn} \neq \delta_{nm}}$.

(2) The δ 's are not algebraically independent:

We have, in fact, for any particular m ,

$$\begin{aligned} & \delta_{m1} + \delta_{m2} + \dots + \delta_{mn} + \dots + \delta_{mt} \\ &= \sum_{j=1}^t \delta_{mj} \\ &= (AmB_1) - \frac{(Am)(B_1)}{N} + (AmB_2) - \frac{(Am)(B_2)}{N} + \dots \\ &\quad + (AmB_t) - \frac{(Am)(B_t)}{N} \\ &= (Am) - \frac{(Am)}{N} \{ (B_1) + (B_2) + \dots + (B_t) \} \\ &= (Am) - \frac{(Am)}{N} \times N \\ &= 0. \end{aligned}$$

OR

$$\begin{aligned} \left[\sum_{j=1}^t \delta_{mj} \right] &= \sum_{j=1}^t \left\{ (AmB_j) - \frac{(Am)(B_j)}{N} \right\} \\ &= \sum_{j=1}^t (AmB_j) - \frac{(Am)}{N} \sum_{j=1}^t (B_j) \\ &= (Am) - \frac{(Am)}{N} \times N \\ &= 0. \end{aligned}$$

\therefore For a particular m , $\sum_{i=1}^s \delta_{im} = 0$.

Now, there are s δ -quantities. In virtue of the relationship we have just proved, for any particular m only $(t-1)$ of the t δ -quantities δ_{mn} are independent. Similarly, for any m only $(s-1)$ are independent. Hence, the total number of independent δ 's is $(s-1)(t-1)$.

SOME MEASURES OF ASSOCIATION : → (C.V.)

(1) Mean Square Contingency : → We may define a measure of association in terms of so called 'square contingency'

$$\chi^2 = \sum_{m=1}^s \sum_{n=1}^t \left(\frac{\delta_{mn}}{(A_m B_n)_0} \right)$$

and call χ^2 the "square contingency".

We then define,

$$\phi^u = \frac{\chi^2}{N}$$

and call ϕ^u the "mean-square contingency".

Clearly, χ^2 and ϕ^u , being the sums of squares, can't be negative. Then vanish if, and only if, every δ -number vanished, in which case

A and B are independent.

& ϕ^u is not suitable for comparative study.

$$\begin{aligned} \chi^2 &= \sum_{m=1}^s \sum_{n=1}^t \frac{\delta_{mn}}{(A_m B_n)_0} \\ &= N \left\{ \sum_{m=1}^s \sum_{n=1}^t \frac{(A_m B_n)^2}{(A_m)(B_n)} - 1 \right\} \\ &= N \left\{ \sum_{m=1}^s \frac{(A_m B_m)^2}{(A_m)(B_m)} - 1 \right\} \quad [\text{where } s=t] \\ &= N(t-1) \\ \Leftrightarrow \phi^u &= t-1. \text{ Hence the limits of } \phi^u \text{ vary in different systems.} \end{aligned}$$

~~Non-parametric coefficient of mean-square contingency~~

Properties : → (a) ϕ^u is non-negative since it is the sum of squares and ϕ^u can take any non-negative any non-negative value, i.e. $0 \leq \phi^u \leq \infty$.

(b) $\phi^u = 0$ iff $\delta_{mn} = 0 \forall m, n$.

iff A and B are independent.

Limitations : → (a) Since χ^2 can take any non-negative value then the limits of ϕ^u vary in different situations. Hence ϕ^u is not suitable to form a coefficient of association.

(b) In ext table, in case of complete association so that $(A_m) = (B_n) = (A_m B_n) \forall m \& (A_m B_n) = 0 \forall m \neq n$, i.e. when only the leading diagonal frequencies are non-zero.

(15)

$$\begin{aligned}
 \underline{\text{Why?}} \quad & \sum_{m=1}^s \sum_{n=1}^t \frac{s v_{mn}}{(A_m B_n)_0} \\
 & = \sum_{m=1}^s \sum_{n=1}^t \frac{N}{(A_m)(B_n)} \left[(A_m B_n) - \frac{(A_m)(B_n)}{N} \right] \\
 & = N \sum_{m=1}^s \sum_{n=1}^t \left\{ (A_m B_n)^2 - \frac{2(A_m B_n)(A_m)(B_n)}{N} + \frac{(A_m)^2 (B_n)^2}{N^2} \right\} / (A_m)(B_n) \\
 & \approx N \sum_{m=1}^s \sum_{n=1}^t \left\{ (A_m B_n)^2 + \frac{(A_m)(B_n)}{N^2} ((A_m)(B_n) - 2(A_m B_n) N) \right\} / (A_m)(B_n) \\
 & = N \left[\sum_{m=1}^s \sum_{n=1}^t \left\{ \frac{(A_m B_n)^2}{(A_m)(B_n)} \right\} + \frac{1}{N^2} \sum_{m=1}^s \sum_{n=1}^t \left\{ (A_m)(B_n) - 2(A_m B_n) N \right\} \right] \\
 & = N \left[\sum_{m=1}^s \sum_{n=1}^t \frac{(A_m B_n)^2}{(A_m)(B_n)} - 1 \right]
 \end{aligned}$$

(c.v.)

(2) Karl Pearson's coefficient of mean-square contingency:

The quantity ϕ^2 is not quite suitable in itself to form a coefficient, because its limits vary in different cases. The upper-limit is infinite as N increases. Pearson, therefore, proposed the coefficient C , defined by

$$C = \sqrt{\frac{\chi^2}{N + \chi^2}} = \sqrt{\frac{\phi^2}{1 + \phi^2}}$$

This is called the "coefficient of mean-square contingency". In general, no sign should be attached to the root, for the coefficient merely shows whether two characters are or are not independent; but in certain cases a conventional sign may be used. As for example — slight pigmentation of eyes and hair appear to go together, and the contingency may be regarded as positive. If slight pigmentation of eyes had been associated with marked pigmentation of hair, the contingency might have been regarded as negative.

(2) Properties:-

(1) C vanishes iff $\delta_{mn} = 0 \forall (m, n)$
i.e. $C = 0$ iff A and B are independent of each other.

(2) The coefficient has one serious disadvantages.

c.v. $\therefore 0 \leq C < 1$ since $0 \leq X^w < N + X^w$

Limitations:-

(1) The coefficient C has one serious drawback that it never reaches the value 1 even in the case of perfect association.

(2) In $t \times t$ table, in case of complete association,
 $X^w = N(t-1)$ and

$$C = \sqrt{\frac{N(t-1)}{N + N(t-1)}} = \sqrt{\frac{t-1}{t}}$$

Note that, no greater association than the case.

$$(A_m B_m) = (A_m) = (B_m) \forall m, \text{ and}$$

$$(A_m B_n) = 0 \forall m \neq n$$

can be imagined in $t \times t$ table. Hence, $X^w \leq N(t-1)$ and
 $C \leq \sqrt{\frac{t-1}{t}}$. The upper limits of C vary for different systems
and hence it is not suitable for comparative study.

c.v. (3) Tschuprow's Coefficient: ~ Even in case of complete association the value of C is not unity and its maximum depends on the number of rows and columns of the table.
To remedy the defect to which we have referred above, Tschuprow proposed an alternative coefficient T ;

defined by,

$$T^w = \left\{ \frac{X^w}{N \sqrt{(s-1)(t-1)}} \right\}$$

$$= \frac{\phi^w}{\sqrt{(s-1)(t-1)}}$$

$T = 1$, in case of complete association in a $t \times t$ table but it will not be true in case of an $s \times t$ classification where $s \neq t$.

The coefficient varies between 0 and 1.

(4) ~~Measures~~

(17)

Cramen's V^2 :

In an $s \times t$ table, if $s > t$ and $(A_m B_m) = (A_m) (B_m)$,
 $m = 1(1)t$; then $\chi^2 = N(t-1)$. But if no broad

$s < t$ & $(A_m B_m) = (A_m) (B_m)$, $m = 1(1)t$; then
 $\chi^2 = N(s-1)$.

Hence for any $s \times t$ table, $\chi^2 \leq N \min\{s-1, t-1\}$.

As Cramers pointed out, we may avoid these difficulties by defining a new measure:

$$V^2 = \frac{\chi^2}{N \min\{s-1, t-1\}}$$

Clearly V^2 lies between 0 and 1.

When, $s=t=2$, $V^2 = \frac{\chi^2}{N} = \phi^2$

Evidently, $V^2 = T^2$ when the table is square, otherwise V^2 exceeds T^2 . Although the difference won't be very large unless s and t varies widely.

We also see that —

$$\frac{C^2}{T^2} = \frac{\sqrt{(s-1)(t-1)}}{1 + \phi^2}$$

Remark: These measures of association, discussed above, are applicable to both nominal and ordinal data.

ORDINAL TRENDS : Let us consider two jointly distributed attributes A and B. Further assume

that the level of each can be arranged according to their degree of possession. Here it may happen that — as the responses on the attribute A increase towards its higher level, responses on B also increase towards its higher level or the responses on B decrease towards its lower level. Such type of association is often referred to as a monotone trend association.

Concordance & Discordance :

A pair is said to be concordant pairs if the subject ranked higher on A also ranked higher on B. A pair is said to be discordant pairs if the subject ranking higher on A, ranks lower on B. A pair is said to be tied, if the subjects have same classification on A and / or on B.

Ordinal measures of association:

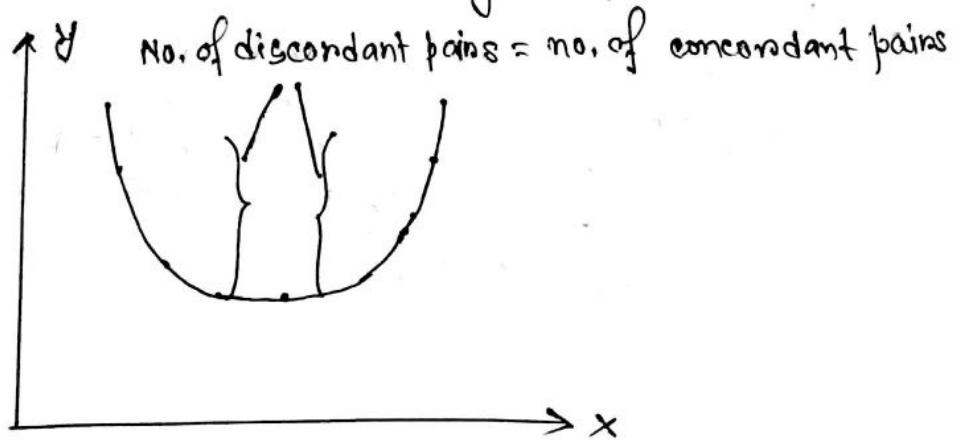
Here we shall consider some measures based on the no. of concordance and the no. of discordance pairs.

(1) Goodman-Kruskal Gamma: Let C & D be respectively the no. of concordance pairs and the no. of discordance pairs. If the pairs are united on both variables, then $\frac{C}{C+D}$ is the proportion of concordance pairs and $\frac{D}{C+D}$ is the proportion of discordance pairs. The difference $(\frac{C}{C+D} - \frac{D}{C+D})$ = $(\frac{C-D}{C+D})$ is termed as Goodman-Kruskal Gamma (γ). i.e. $\gamma = \frac{C-D}{C+D}$.

If $\gamma \geq 0$, i.e. $\frac{C}{C+D} \geq \frac{D}{C+D}$, then the majority of pairs are in concordance (or discordance), i.e. characters A and B have the +ve (or, -ve) association.

Properties:

- The Goodman-Kruskal Gamma (γ) treats the characters symmetrically — any one of them can be taken as response variable.
- Clearly, $|\gamma| \leq 1$, i.e. $-1 \leq \gamma \leq 1$, since $C \geq 0$ and $D \geq 0$.
- Under independence γ vanishes but the converse is not true. i.e. Independence $\Rightarrow \gamma = 0$ but the converse is not necessarily true. For a U-shaped joint-distrn., one can have $(C = D)$, i.e. $\gamma = 0$ and the characters are not independent (pic: (c))
- $\gamma = \pm 1$ iff the characters are monotonically related.
 $\gamma = 1 \Leftrightarrow D = 0$ iff the characters have monotonic increasing relation.
 $\gamma = -1 \Leftrightarrow C = 0$ iff the characters have monotonic decreasing relation.



pic: (c)

REMARK: Note that,

$$C = 2 \sum_i \sum_j n_{ij} \left(\sum_{k>i} \sum_{l>j} n_{ikl} \right)$$

$$D = 2 \sum_i \sum_j n_{ij} \left(\sum_{k>i} \sum_{l>j} n_{lki} \right)$$

Let, $T_A = \text{No. of tied pairs on the characters A}$

$$= \binom{n_{10}}{2} + \binom{n_{20}}{2} + \dots + \binom{n_{80}}{2}$$

$$= \sum_{i=1}^8 \frac{n_{10}(n_{10}-1)}{2}, \text{ since } n_{10} \text{ members have}$$

the same characters A.

Similarly, $T_B = \text{No. of tied pairs on the characters B}$,

$$= \binom{n_{01}}{2} + \binom{n_{02}}{2} + \binom{n_{03}}{2} + \dots + \binom{n_{0t}}{2}$$

$$= \sum_{j=1}^t \frac{n_{0j}(n_{0j}-1)}{2}$$

Now, $T_{AB} = \text{No. of tied pairs on both A and B}$

$$= \sum_{i=1}^8 \sum_{j=1}^t \binom{n_{ij}}{2}$$

$$= \sum_{i=1}^8 \sum_{j=1}^t \frac{n_{ij}(n_{ij}-1)}{2}$$

$$\text{Note that } \binom{N}{2} = C + D + T_A + T_B - T_{AB}$$

(2) Kendall's T : Let (x_i, y_i) and (x_j, y_j) be two pairs of observation on two characters for a pair (i, j) , $i < j$ of individuals.

Define, $a_{ij} = \text{sign}(x_i - x_j)$

$$= \begin{cases} 1 & \text{if } x_i > x_j \\ 0 & \text{if } x_i = x_j \\ -1 & \text{if } x_i < x_j \end{cases}$$

and, $b_{ij} = \text{sign}(y_i - y_j)$

Then, our measure of rank correlation will be based on the sum,

$$S = \sum_{i < j} a_{ij} b_{ij} - \binom{n}{2} = \sum_{i < j} a_{ij} b_{ij} - \frac{n(n-1)}{2}$$

$= C - D$, which is the total score.

Hence, we have two choices to make a measure of association, if we wish to standardize S to lie in the range [-1, 1] and attain ± 1 in the extreme cases of complete disassociation and complete association.

(a) Kendall's T_a : — If there are no ties, no a_{ij} or b_{ij} are zero. In that case, S could vary in between $-(\frac{N}{2})$ & $(\frac{N}{2})$. i.e. $-(\frac{N}{2}) \leq S \leq (\frac{N}{2})$. Hence, we define,

$$\text{Kendall's } T_a = \frac{S}{(\frac{N}{2})}, \text{ as a coefficient of association.}$$

Clearly, whether T_a attains the end points or not, completely depends on the no. of zero scores.

If some scores a_{ij} or b_{ij} are zero, the T_a can no longer attain ± 1 .

$$S = \sum_{i < j} \sum a_{ij} b_{ij} = C - D$$

If there are no ties, then $S = \pm (\frac{N}{2})$ if all pairs are concordant or discordant.

If there are some ties, then —

$$-(\frac{N}{2}) < S < (\frac{N}{2}) \text{ &}$$

Consequently, T_a can't attain ± 1 .

Note that, $T_a = \frac{S}{(\frac{N}{2})} = \frac{C-D}{(\frac{N}{2})}$, which is same as the Kendall's T_a .

(b) Kendall's T_b : — The correlation coefficient between two set of scores (a_{ij}) and (b_{ij}) is defined as —

$$T_b = \frac{\sum_{i < j} \sum a_{ij} b_{ij}}{\sqrt{\sum_{i < j} \sum a_{ij}^2} \sqrt{\sum_{i < j} \sum b_{ij}^2}}$$

In the presence of tied pairs, the measure T_b will be of the form, $T_b = \frac{C-D}{\sqrt{(\frac{N}{2})-T_A} \sqrt{(\frac{N}{2})-T_B}}$

Now, $\sum_{i < j} \sum a_{ij}^2$ = No. of pairs (i, j) , $i < j$ for which $x_i > x_j$ or $x_i < x_j$
 $= (\frac{N}{2}) - \text{No. of tied pairs on A}$
 $= (\frac{N}{2}) - T_A$, where $T_A = \sum_{i=1}^N \frac{n_{i0}(n_{i0}-1)}{2}$

& similarly, $\sum_{i < j} \sum b_{ij}^2 = (\frac{N}{2}) - T_B$.

If there are no ties, then $T_a = T_b$, but if there are some ties then $T_b > T_a$.

Remark: → The Goodman-Kruskal γ can also be expressed as —

$$\gamma = \frac{C - D}{C + D} = \frac{C - D}{\binom{N}{2} - T_A - T_B + T_{AB}}$$

to note

It is important that γ, T_A, T_B all attain ± 1 if all the observations are in the cells along the longest diagonal of the table.

(3) Somers' d : — Somers proposed a measure of association for ordinal data, define,

$$d = \frac{C - D}{\binom{N}{2} - T_A} \text{ or } \frac{C - D}{\binom{N}{2} - T_B} \rightarrow \text{as a coefficient of association.}$$

In case of complete association, $(A_m B_m) = (A_m) = (B_m)$, $m = 1(i) \min(s, t)$, and all other cell frequencies are zero, i.e. all the observations lie in a longest diagonal of a table. Then

In that case,

$$S = \sum_{i < j} \sum a_{ij} b_{ij} = \left\{ \begin{array}{l} \text{No. of pairs } (i, j), i < j \text{ for which} \\ a_i < a_j \text{ (or } y_i < y_j \text{ or } x_i > x_j \text{)} \end{array} \right\}$$

$$\text{freedom} = \left\{ \begin{array}{l} \text{No. of pairs united on A (or B)} \\ = \left\{ \binom{N}{2} - T_A \right\} \text{ or } \left\{ \binom{N}{2} - T_B \right\} \end{array} \right\}.$$

Hence, $d = 1$.

But if all the observations lie on the longest diagonal of a table such that

$$(A_m B_{s-m+1}) = (A_m) = (B_{s-m+1}) \quad \forall m = 1(i) \min(s, t)$$

and all other frequencies are zero. (Assuming that $s < t$),

$$\text{Then, } S = \sum_{i < j} \sum a_{ij} b_{ij} = - \left\{ \begin{array}{l} \text{No. of pairs } (i, j), i < j \text{ for which} \\ x_i > x_j \text{ (or } y_i > y_j \text{)} \end{array} \right\}$$

$$= - \left\{ \binom{N}{2} - T_B \right\} \text{ on } - \left\{ \binom{N}{2} - T_A \right\}$$

and then $d = -1$.

Remark: → In the above situations, Kendall's T_B behaves similarly.



Logistic Regression

Consider a clinical trial

where patients are given a treatment and response ($=1$ if success; $=0$ if failure) is observed. In addition their age is also reported. Let y be the response variable and x be the age. Then we assumed

$$P(y=1) = \pi(x)$$

$$\pi(x) = \frac{e^{ax+bx}}{1+e^{ax+bx}} ; a, b \text{ are unknown.}$$

Consider m points of data $(x_i, y_i), i=1(1)m$. Then the corresponding regression is termed as "Logistic regression".

$$a+bx = \log \frac{\pi(x)}{1-\pi(x)} = \text{log odds ratio representation.}$$

Here, we can't use linear regression model i.e. $y=a+bx$, because for the unknown a and b may not take the value 0 or 1.

Ref:- Agresti :-

Let Y denote a response variable that can assume only two values, say 0 and 1. Denote the expected value of Y by $E(Y) = P(Y=1) = \pi$

and suppose that we want to model $\pi(x)$ on the values of explanatory variables $x = (x_1, x_2, \dots, x_k)$. The standard regression model has the form

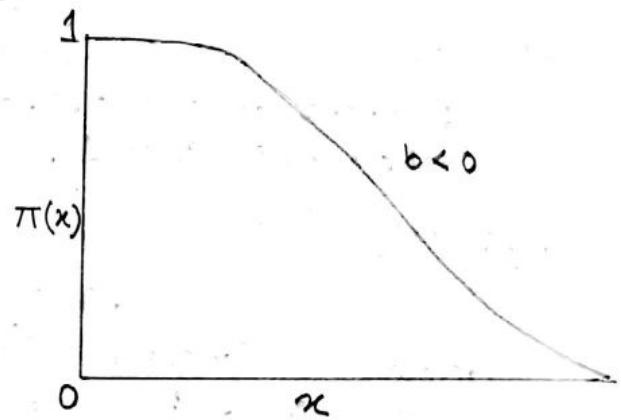
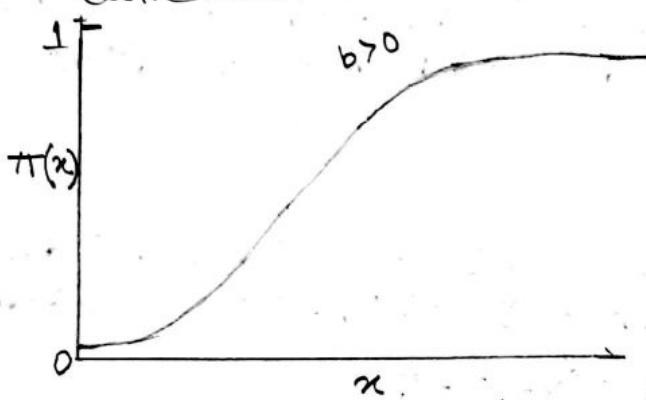
$$\pi(x) = a + b_1 x_1 + \dots + b_k x_k$$

There are several difficulties with using ordinary least square principle to fit a model of standard regression conditions that make least squares estimates optimal are not these. For instance, the variance of Y is $\pi(x)(1-\pi(x))$, which is not constant over the range of values of the explanatory variables. Standard distributional statements for estimations do not apply, since Y is dichotomous rather than normally distributed.

$$\sigma^2 = \left[\frac{\pi}{(1-\pi)} \right] \text{Var}(Y)$$

A weighted least squares approach can be used to obtain more efficient estimates of the regression parameters in this model. The model itself is likely to be inaccurate in certain regions, however, if some x_i 's are quantitative. They follow because the model predicts the impossible value $\pi < 0$ and $\pi > 1$ for sufficiently large or sufficiently small values of x_i . For a dichotomous response, $E(Y)$ can't be linearly related to x_i over an unbounded range of x_i values.

Logistic function : — Because of the factors just discussed, it is often more appropriate to use a model that allows a curvilinear relationship between $E(Y)$ and each quantitative x_i . If we expect monotonic relationship then the regression curve will be



Here $\pi(x) = \frac{e^{ax+bx}}{1+e^{ax+bx}}$ called the logistic function. This function is monotonic with $\pi(x) \downarrow 0$ or $\pi(x) \uparrow 1$ as $x \uparrow \infty$ depending on whether $b < 0$ or $b > 0$. Takes the value $\pi(x) = \frac{1}{2}$ at $x = -\frac{a}{b}$ and the curve has a steeper rate of increase around that value as $|b|$ increases. When $b > 0$, it is the distribution function of the logistic random variable having mean $-\frac{a}{b}$ and s.d. $\frac{\pi}{\sqrt{3}b}$.

Hence the odds of making response 1 instead of response 0 is

$$\frac{\pi(x)}{1-\pi(x)} = e^{ax+bx}$$

The odds increases multiplicatively by e^b for every unit of x . The log odds has the simple linear relationship

$$\log \left[\frac{\pi(x)}{1-\pi(x)} \right] = ax+bx$$

The model of the log odds is called logistic regression model. In this relationship the logit transformation yields a linear relationship for the logit model. When there are several explanatory variables, the logit model generalizes to

$$\log \left[\frac{\pi(x)}{1-\pi(x)} \right] = a + b_1 x_1 + \dots + b_k x_k.$$

Logit model: → In applied mathematics & statistics, the logit of a number p between 0 and 1 is

$$\text{logit}(p) = \log \left(\frac{p}{1-p} \right) = \log p - \log(1-p)$$

The logit function is the inverse of the sigmoid or logistic function. If p is a probability then $p/(1-p)$ is the corresponding odds and the logit of the probability is the logarithm of the odds.

Q. Why do we use logistic regression rather than ordinary linear regression?

Ans:—

→ If we use linear regression, the predicted values will become greater than one and less than zero if you move far enough on the x-axis. Such values are theoretically inadmissible.

→ One of the assumption of regression is that variance of y is constant across values of x . This cannot be the case with a binary variable because the variance is $p(1-p)$. So, if p approaches to 1 or 0, variance approaches to zero.

→ The significance of testing of the weights b rest upon the assumption that errors of prediction ($y - \hat{y}$) are normally distributed. Because y only takes values 0 and 1, this assumption is hard to justify. Therefore, we can't use linear regression.

Relative Risk: — A value $\pi_1 - \pi_2$ of fixed size may have greater importance, when both π_i are close to zero or 1 than when they are not. For a study comparing two treatments on the proportion of subjects who die, the difference between 0.010 and 0.001 and 0.410 and 0.401 are both 0.009. In such cases, the ratio of proportions is also informative. The relative risk is defined to be the ratio $\frac{\pi_1}{\pi_2}$. It can be any non-negative real number. Relative risk of one corresponds to independence, e.g. for the proportions just given the relative risks are $\frac{0.010}{0.001} = 10$ and $\frac{0.410}{0.401} = 1.02$.

Comparing the rows on the 2nd response (failure) category gives a different relative risk, $\frac{1-\pi_1}{1-\pi_2}$.

(c) Odds Ratio: — For a probability π of a success, the odds are defined to be $\Omega = \frac{\pi}{1-\pi}$. The odds are non-negative with $\Omega > 1$ when a success is more likely than a failure, e.g.

$\pi = 0.75 \Rightarrow \Omega = 3$, a success is thrice as likely as a failure.

Inversely, if $\Omega = \frac{\pi}{1-\pi}$,

Refer again to 2x2 table, within row i , the odds of success instead of failure are $\Omega_i = \frac{\pi_{i1}}{1-\pi_{i1}}$.

The ratio of the odds Ω_1 and Ω_2 in the two rows,

$\Theta = \frac{\Omega_1}{\Omega_2} = \frac{\pi_1 / (1-\pi_1)}{\pi_2 / (1-\pi_2)}$; is called the odds ratio.

For joint distn. with cell probabilities $\{\pi_{ij}\}$, the

equivalent defn. of Ω_i could be

$\Omega_i = \frac{\pi_{i1}}{\pi_{i2}}$. Hence,

$$\Theta = \frac{\pi_{11} / \pi_{12}}{\pi_{21} / \pi_{22}} = \frac{\pi_{11} \pi_{22}}{\pi_{12} \pi_{21}}$$

Hence, odds ratio is also called cross-product ratio.

Properties of Odds Ratio: — The odds ratio can be equal to any non-negative number, the condition $\Omega_1 = \Omega_2$ and hence (when all cell probabilities are +ve), $\Omega = 1$ corresponds to independence of X and Y . When $1 < \Omega < \infty$ subjects in row 1 are more likely to have a success than subjects in row 2; i.e. $\pi_{11} > \pi_{21}$. For instance, when $\Omega = 4$, the odds of success in row 1 are 4 times the odds in row 2, this does not mean that $\pi_{11} = 4\pi_{21}$, i.e. relative risk of 4. When $0 < \Omega < 1$, $\pi_{11} < \pi_{21}$, one of the cell probabilities vanishes, Ω becomes 0 or ∞ .

Values of Ω further from 1 in a given direction represents stronger association. Two values represent the same association, but in opposite directions, when one is the inverse of the others.

c.v.

Relationship between Odds ratio & Relative Risk: —

$$\text{Odds ratio} = \text{relative risk} \times \frac{1 - \pi_{22}}{1 - \pi_{11}}$$

Their magnitude are similar whenever the probability π_{ii} of the outcome of interest is close to zero for both groups thus when each π_{ii} is small, the odds ratio provides a rough indication of the relative risk when it is not directly estimable, such as in case of controlled study.

We can compute the odds ratio, however, since it is determined by the conditional distributions in either direction. When the probability of the outcome of interest is very small, the population odds ratio and relative risk take similar values.

Logistic Regression:

Let Y denote a binary response variable.

variable, for instance Y might indicate the choice of car (new or used) or the diagnosis of cancer (present or absent). Each observation has one of 2 outcomes denoted by 0 and 1. Now for a response variable Y and an explanatory variable X , consider the data.

Setting of x	values of Y	Totals
x_1	$y_{11} y_{12} \dots y_{1N_1}$	$\sum_{j=1}^{N_1} y_{1j} = n_1$
x_2	$y_{21} y_{22} \dots y_{2N_2}$	$\sum_{j=1}^{N_2} y_{2j} = n_2$
\vdots		
x_K	$y_{K1} y_{K2} \dots y_{KN_K}$	$\sum_{j=1}^{N_K} y_{Kj} = n_K$

Here the true regression of y on x is given by the array mean:

$$\pi(x_i) = \bar{y}_i = \frac{\sum_{j=1}^{N_i} y_{ij}}{N_i} = \frac{n_i}{N_i}$$

Since n_i is the number of y_{ij} which takes the value 1 when $x=x_i$ and $0 \leq n_i \leq N_i$. Hence, $0 \leq \bar{y}_i \leq 1$ and

$$0 < \frac{\bar{y}_i}{1-\bar{y}_i} < \infty \text{, then, } -\infty < \log_e \left(\frac{\bar{y}_i}{1-\bar{y}_i} \right) < \infty.$$

i.e. $\log_e \left(\frac{\bar{y}_i}{1-\bar{y}_i} \right)$ can be any real numbers, the real nos. are also the range of any linear predictor such as $(\alpha + \beta x)$. We can assume the regression model as $\log_e \left(\frac{\bar{y}_i}{1-\bar{y}_i} \right) = \alpha + \beta x_i$:

$$\Leftrightarrow \bar{y}_i = \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}}$$

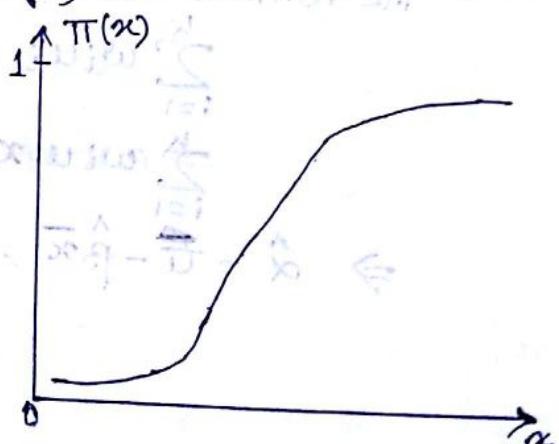
Usually binary data result from a non-linear relationship between $\pi(x_i) = \bar{y}_i$ and x_i . A fixed change in x has less impact on $\pi(x)$ when $\pi(x)$ is near 0 or 1 than when $\pi(x)$ is near 0.5.

In practice, non linear relationship between $\pi(x)$ and x is monotonic. Most important curve with the above shape is the logistic curve.

An appropriate regression model is

$$\pi(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}},$$

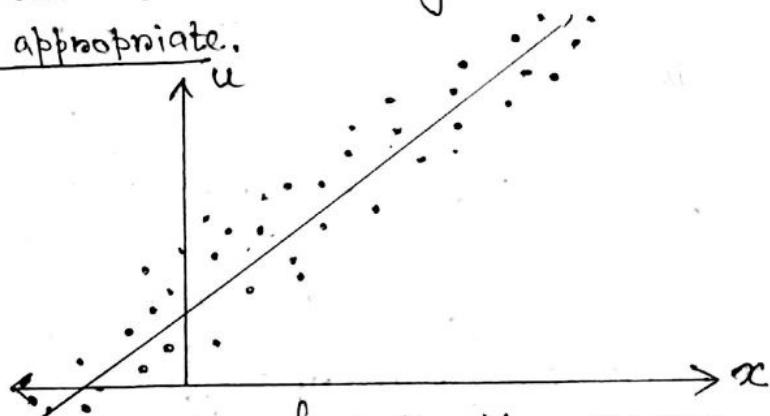
which is called logistic regression model.



looking at the data & fitting of logistic regression equation

Before fitting logistic model, look at the data to check whether the logistic regression is appropriate or not. Since Y takes any of the values 0 and 1, it is difficult to check by plotting Y -values for different i , it is difficult to check by plotting y -values for different i .

If the quantities $u_i = \log\left(\frac{\bar{y}_i}{1-\bar{y}_i}\right) = \log\left(\frac{n_i}{N_i-n_i}\right)$ are plotted against x_i , $i=1(1)k$. and the points (x_i, u_i) are near about a line. Then, by definition, the logistic regression is appropriate.



- Fitting: Consider the predicting formula $U_x = \alpha + \beta x$. When $x = x_i$, then the predicting value of $u_i = \log\left(\frac{\bar{y}_i}{1-\bar{y}_i}\right)$ is $U_{x_i} = \alpha + \beta x_i$ and the corresponding error $e_i = u_i - U_{x_i} = (u_i - \alpha - \beta x_i)$. To determine α and β , we shall use method of weighted least squares and we shall minimise.

$S^2 = \sum_{i=1}^k w_i e_i^2 = \sum_{i=1}^k w_i (u_i - \alpha - \beta x_i)^2$, where the weights $w_i = N_i \bar{y}_i (1 - \bar{y}_i) = N_i \pi(x_i) (1 - \pi(x_i))$, with respect to α and β .

The normal equations are:

$$\sum_{i=1}^k w_i u_i = \alpha \sum_{i=1}^k w_i + \beta \sum_{i=1}^k w_i x_i$$

$$\sum_{i=1}^k w_i u_i x_i = \alpha \sum_{i=1}^k w_i x_i + \beta \sum_{i=1}^k w_i x_i^2$$

$$\Rightarrow \hat{\alpha} = \bar{u} - \hat{\beta} \bar{x}, \quad \hat{\beta} = \frac{\sum_{i=1}^k w_i u_i x_i - N \bar{u} \bar{x}}{\sum_{i=1}^k w_i x_i^2 - N \bar{x}^2}, \text{ where}$$

$$\bar{u} = \frac{\sum w_i u_i}{\sum w_i} \quad \text{and} \quad \bar{x} = \frac{\sum w_i x_i}{\sum w_i}.$$

Question: — What is the rationale behind the Yule's coefficient of association?

Soln. → δ is a measure of independence of A and B. To get a normalized measure, we divide $N\delta$, the difference between $(AB)(\alpha\beta)$ & $(A\beta)(\alpha B)$ by the sum $(AB)(\alpha\beta) + (A\beta)(\alpha B)$ and the resulting coefficient is $Q = \frac{(AB)(\alpha\beta) - (A\beta)(\alpha B)}{(AB)(\alpha\beta) + (A\beta)(\alpha B)}$, which is suitable for comparison of two or more data sets.

Question: — What is the rationale behind Yule's coefficient of colligation?

Soln. → Empirical study reveals that the coefficient Q overestimate the extent of association. Hence, Yule provides a measure by considering lower orders of the cell frequencies as γ , a coefficient of colligation. Note that $Q = \frac{2\gamma}{1+\gamma} \geq \gamma$.

Question: — Let $x_i = \begin{cases} 1, & \text{i-th individual possess A} \\ 0, & \text{i-th individual possess } \alpha \end{cases}$

$y_i = \begin{cases} 1, & \text{i-th individual possess B} \\ 0, & \text{i-th individual possess } \beta \end{cases}$

Find the correlation coefficient between x and y based on the data $\{(x_i, y_i) : i=1(1)n\}$.

$$\text{Soln.} \rightarrow \sum_{i=1}^n x_i y_i = n\bar{x}\bar{y}$$

$$r_{xy} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2\right)^{1/2} \left(\sum_{i=1}^n y_i^2 - n\bar{y}^2\right)^{1/2}}$$

$$f_{11} = n \cdot \frac{f_{10}}{n} \cdot \frac{f_{01}}{n}$$

$$= \frac{\left\{ f_{10} - \frac{f_{10}}{n} \right\}^{1/2} \left\{ f_{01} - \frac{f_{01}}{n} \right\}^{1/2}}{\left(f_{10} - \frac{f_{10}}{n} \right)^{1/2} \left(f_{01} - \frac{f_{01}}{n} \right)^{1/2}}$$

$$= \frac{(f_{11} f_{22} - f_{12} f_{21})/n}{\left(\frac{f_{10} f_{20}}{n} \right)^{1/2} \left(\frac{f_{01} f_{02}}{n} \right)^{1/2}}$$

$$= \frac{f_{11} f_{22} - f_{12} f_{21}}{\sqrt{f_{10} f_{20} f_{01} f_{02}}}$$

$$= \frac{n\delta}{\sqrt{f_{10} f_{20} f_{01} f_{02}}}$$

$$= \frac{n\delta}{\sqrt{f_{10} f_{20} f_{01} f_{02}}} \quad | \quad \begin{array}{c} x \\ y \end{array} \quad | \quad \begin{array}{cccc} 1 & & & \\ & f_{11} & f_{21} & f_{01} \end{array}$$

$$= \frac{n\delta}{\sqrt{f_{10} f_{20} f_{01} f_{02}}} \quad | \quad \begin{array}{cccc} 0 & & & \\ & f_{12} & f_{22} & f_{02} \end{array}$$

$$= \frac{n\delta}{\sqrt{f_{10} f_{20} f_{01} f_{02}}} \quad | \quad \begin{array}{cc} f_{10} & f_{20} \\ f_{10} & n \end{array}$$

Hence, $r_{xy} = 0$ iff $\delta = 0$ iff A & B are independent.

Problem: For a 2×2 table, Yule introduces

$$\hat{\theta} = \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{11}n_{22} + n_{12}n_{21}}$$

i) S.T. in a 2×2 table, Goodman-Kruskal's $\gamma = \hat{\theta}$.

ii) S.T. the Yule coefficient can be written as

$\hat{\theta} = \frac{\theta - 1}{\theta + 1}$, a monotonic transformation of θ from $[0, \infty)$ to $[-1, 1]$.

Soln. \Rightarrow

i)

$$\gamma = \frac{C - D}{C + D}$$

Hence, C = No. of pairs (i, j) , $i < j$ for which

$\{x_i > x_j, y_i > y_j\}$ or $\{x_i < x_j, y_i < y_j\}$

$$= n_{11} \times n_{22}$$

and D = No. of pairs (i, j) , $i < j$ for which $\{x_i > x_j, y_i < y_j\}$ or $\{x_i < x_j, y_i > y_j\}$

$$= n_{21} \times n_{12}$$

$$\therefore \gamma = \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{11}n_{22} + n_{12}n_{21}} = \hat{\theta}$$

ii) $\hat{\theta}$ = the sample odds ratio on α -measure

$$= \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

$$\therefore \hat{\theta} = \frac{\hat{\theta} - 1}{\hat{\theta} + 1},$$

$$\frac{d\hat{\theta}}{d\hat{\theta}} = \frac{d}{d\hat{\theta}} \left(1 - \frac{2}{\hat{\theta} + 1} \right) = \frac{2}{(\hat{\theta} + 1)^2} > 0$$

$\Rightarrow \hat{\theta}$ is monotonic increasing in $\hat{\theta}$.

As $0 \leq \hat{\theta} < \infty$, $1 \leq \hat{\theta} + 1 < \infty$

$$\Rightarrow 0 < \frac{2}{\hat{\theta} + 1} \leq 2$$

$$\Rightarrow -1 \leq 1 - \frac{2}{\hat{\theta} + 1} \leq 1$$

$$\Rightarrow -1 \leq \hat{\theta} \leq 1.$$

X	Y	y_1	y_2
x_1		n_{11}	n_{12}
x_2		n_{21}	n_{22}

Categorical Data Analysis

Categorical Table: Let us consider two categorical variables A and B, where A takes values on p categories, say, A_1, A_2, \dots, A_p and B takes values on q categories, B_1, B_2, \dots, B_q , on q categories.

Let us define the probability,

$$\pi_{ij} = P(A = A_i \cap B = B_j) \text{ for } i=1(1)p, j=1(1)q.$$

If we construct a table consisting of p rows for the p-categories of A and q columns for the q-categories of B, then we get a $p \times q$ table as follows:

		B						
		B ₁	B ₂	...	B _j	...	B _q	
A ₁	π_{11}	π_{12}	...	π_{1j}	...	π_{1q}	π_{10}	
	π_{21}	π_{22}	...	π_{2j}	...	π_{2q}	π_{20}	
A ₂	π_{31}	π_{32}	...	π_{3j}	...	π_{3q}	π_{30}	
	\vdots	\vdots	...	\vdots	...	\vdots	\vdots	
A _p	π_{p1}	π_{p2}	...	π_{pj}	...	π_{pq}	π_{p0}	
	π_{01}	π_{02}	...	π_{0j}	...	π_{0q}	π_{00}	

This table is called a $p \times q$ or contingency table. Here —

$$\pi_{10} = \sum_{j=1}^q \pi_{1j} \quad \& \quad \pi_{0j} = \sum_{i=1}^p \pi_{ij}$$

are called Marginal Probabilities.

π_{ij} 's are called the Joint probabilities and with their

probabilities we can define the conditional probabilities in this table as follows:

$$\begin{aligned} \pi_{j|i} &= P(B = B_j | A = A_i) \\ &= \frac{P(A = A_i \cap B = B_j)}{P(A = A_i)} \\ &= \frac{\pi_{ij}}{\pi_{i0}} \end{aligned}$$

		B							
		B ₁	B ₂	...	B _j	...	B _q		
A ₁			π_{11}	π_{12}	...	π_{1j}	...	π_{1q}	π_{10}
			π_{21}	π_{22}	...	π_{2j}	...	π_{2q}	π_{20}
		\vdots	\vdots	...	\vdots	...	\vdots	\vdots	\vdots
		π_{p1}	π_{p2}	...	π_{pj}	...	π_{pq}	π_{p0}	
		π_{01}	π_{02}	...	π_{0j}	...	π_{0q}	π_{00}	

Independence of two categorical variables:

Two categorical variables A and B taking values on two categorical scales consisting of the categories A_1, A_2, \dots, A_p and B_1, B_2, \dots, B_q , respectively, are independent if the conditional probability of assuming a value by one variable, say B, remains the same for any value assumed by the other variable.

$$\text{For every } j, \pi_{j|i} = P(B=B_j | A=A_i) = \pi_{0j} \quad \forall i$$

$$\text{i.e. for every } j, \frac{\pi_{ij}}{\pi_{i0}} = \pi_{0j} \quad \forall i$$

$$\Rightarrow \pi_{ij} = \pi_{i0} \cdot \pi_{0j} \quad \forall i, j$$

Difference measures of independence in a 2×2 contingency table:

Let there be two groups of subjects and let each of them be allowed to response on a binary variable. Let the categorical variable corresponding to the group of subjects denoted by X and that corresponding to the response be denoted by Y. Then both X and Y are binary variables.

Response Groups \\\backslash	Yes	No
Grp-I	π_{11}	$1-\pi_{11}$
Grp-II	π_{21}	$1-\pi_{21}$

In the adjacent table, let π_{i1} ($i=1, 2$) be the proportion of 'Yes' response in the i -th group. Note that here the probabilities π_{11} and π_{21} are actually conditional probability given for the group.

Now if the categorical variable X and Y are independent, then $\pi_{11} - \pi_{21} = 0$. Now considering another table —

Response Groups \\\backslash	Yes	No	No. of subjects
Grp-I	x_1	$n_1 - x_1$	n_1
Grp-II	x_2	$n_2 - x_2$	n_2

Now, we can judge independence from the difference of the observed proportion

$$p_1 = \frac{x_1}{n_1} \text{ and } p_2 = \frac{x_2}{n_2}$$

If $p_1 - p_2$ is close to zero,

then we can predict independence, otherwise, association is believed to exist between the categorical variable X and Y.

Limitations: → Limitation of this measure of independence is that when π is close to 0 or 1, the difference may lead to a wrong decision about independence; for instance, let $\pi_1 = 0.010$ and $\pi_2 = 0.001$, then $\pi_1 - \pi_2 = 0.009$, again let $\pi_1^* = 0.410$ and $\pi_2^* = 0.401$, then $\pi_1^* - \pi_2^* = 0.009$, i.e. difference of probabilities are same in both the cases, though $\frac{\pi_1}{\pi_2} = 10$ and $\frac{\pi_1^*}{\pi_2^*} = 1.2$. Limits of this measure are -1 and 1 , $\pi_1 - \pi_2 = +1 \Rightarrow$ extreme +ve association, $\pi_1 - \pi_2 = -1 \Rightarrow$ extreme -ve association.

Example:-

Groups	AGES	YES	NO	TOTAL	
Grp-I	Age > 70	60	0	60 = n_1	$\pi_1 = 1; \pi_2 = 0$ $\therefore \pi_1 - \pi_2 = 1$
Grp-II	Age ≤ 20	0	40	40 = n_2	$\pi_1 = 0; \pi_2 = 1$ $\therefore \pi_1 - \pi_2 = -1$ where, $\pi_1 = \frac{x_1}{n_1}$, $\pi_2 = \frac{x_2}{n_2}$.

ODDS RATIO:- Let us compare two groups of subjects on a binary response variable Y . Let the categorical variable Y be allowed to take values 1 and 2. Then we denote X another categorical variable, also taking values 1 and 2, to denote the two groups of subjects :

X	Response (Y)		1
	1	2	
1	π_1	$1 - \pi_1$	1
2	π_2	$1 - \pi_2$	1

If, we consider the response '1' as "success" with probability π , then the odds of success is defined by $\frac{\pi}{1-\pi}$. In the present case the odds of success for the 1st group is

$$\text{Odds}_1 = \frac{\pi_1}{1 - \pi_1} \quad \text{and that for the 2nd group is}$$

$$\text{Odds}_2 = \frac{\pi_2}{1 - \pi_2}.$$

Then, we define the Odds Ratio of the two groups by

$$\theta = \frac{\text{Odd } S_1}{\text{Odd } S_2} = \frac{\pi_{11}/(1-\pi_{11})}{\pi_{12}/(1-\pi_{12})} = \frac{\pi_{11}(1-\pi_{12})}{\pi_{12}(1-\pi_{11})}$$

when two categorical variables X and Y are independent, $\pi_{11} = \pi_{12}$, then the odds ratio (θ) = 1.
 When $\theta > 1$, we say that X and Y are +vely associated.
 & $\theta < 1$, we say that X and Y are -vely associated.

PROPERTIES OF ODDS RATIO:-

- 1) Odds Ratio can take any non-negative value.
 - 2) $\theta = 1 \Rightarrow$ Independence.
 $\theta > 1 \Rightarrow$ positive association.
 $\theta < 1 \Rightarrow$ negative association.
 - 3) If θ and θ^* are two odds ratios such that
 $\theta = \frac{1}{\theta^*}$; then we can calculate that
 degree of association is the same
 in both the cases, though the
 direction of association is
 opposite.
 - 4) For a table with joint probabilities, the odds ratio is given by
- $$\theta = \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}}$$

Proof:- In the table, we can write the odds ratio in terms of conditional probabilities

$$\pi_{111} = \frac{\pi_{11}}{\pi_{11} + \pi_{12}} \text{ and } \pi_{112} = \frac{\pi_{21}}{\pi_{21} + \pi_{22}}$$

$$\begin{aligned} \text{Then, } \theta &= \frac{\pi_{111}/(1-\pi_{111})}{\pi_{112}/(1-\pi_{112})} = \frac{\pi_{111}(1-\pi_{112})}{\pi_{112}(1-\pi_{111})} \\ &= \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}}. \end{aligned}$$

\Rightarrow If any rows or any column of the table is interchanged, then the odds ratio is reversed.

Proof:-

	1	2
1	π_1	$1-\pi_1$
2	π_2	$1-\pi_2$

$$\text{Odds ratio is } \theta = \frac{\pi_1(1-\pi_2)}{\pi_2(1-\pi_1)}$$

Now, if the 1st and 2nd rows of the table are interchanged to get the new table.

	1	2
1	π_2	$1-\pi_2$
2	π_1	$1-\pi_1$

Then the odds ratio is

$$\theta^* = \frac{\pi_2(1-\pi_1)}{\pi_1(1-\pi_2)} = \frac{1}{\theta}.$$

6. If the orientation of the table is reversed then the odds ratio does not change.

Proof:-

	1	2
1	π_1	$1-\pi_1$
2	π_2	$1-\pi_2$

Hence the odds ratio

$$\theta = \frac{\pi_1(1-\pi_2)}{\pi_2(1-\pi_1)}$$

Now, if the orientation of this table is reversed to get the new table,

	1	2
1	π_1	$1-\pi_1$
2	π_2	$1-\pi_2$

Here, the odds ratio is

$$\theta^* = \frac{\pi_1/\pi_2}{(1-\pi_1)/(1-\pi_2)} = \frac{\pi_1(1-\pi_2)}{\pi_2(1-\pi_1)}$$

	1	2
1	π_1	π_2
2	$1-\pi_1$	$1-\pi_2$

$$= \frac{\pi_1(1-\pi_2)}{\pi_2(1-\pi_1)}$$

	1	2
1	π_1	$1-\pi_1$
2	π_2	$1-\pi_2$

$$= \theta$$

c.v. PROBLEM : → Show that sample odds ratio does not change even when both cell counts within any row are multiplied by a non-zero constant or when both cell counts within any column are multiplied by a non-zero constant. Discuss the implication of the above result using a real life example.

ANS: →

	1	2	
1	π_1	$1-\pi_1$	
2	π_2	$1-\pi_2$	

	1	2	
1	$k\pi_1$	$k(1-\pi_1)$	
2	π_2	$1-\pi_2$	

Odds ratio is

$$\Theta = \frac{\pi_1(1-\pi_2)}{\pi_2(1-\pi_1)}$$

Now, if the both cell counts within any rows are multiplied by a constant (nonzero), say, k , we get the new table.

then the odds ratio is

$$\Theta^* = \frac{k\pi_1(1-\pi_2)}{\pi_2(k(1-\pi_1))} \\ = \Theta$$

So, odds ratio does not change.

An implication of the multiplicative invariance property is that the sample odds ratio estimates the same characteristic (Θ) even when we select disproportionately large or small samples from marginal categories of a variable. For instance, suppose a study investigates the association between vaccination and catching a certain strain of flu. For a retrospective design, the sample odds ratio estimates the same characteristic whether we randomly sample (1) 100 people who got the flu and 100 people who didn't, (2) 150 people who got the flu and 50 people who didn't, in each case classifying subjects on whether they took the vaccine. In fact, the odds ratio is equally valid for retrospective, prospective, or cross-sectional sampling designs. We would estimate the same characteristic if (3) we randomly sample 100 people who took the vaccine and 100 people who didn't, and then classify them on whether they got the flu, or (4) we randomly sample 200 people and classify them on whether they took the vaccine and whether they got the flu.

— X —