

ST. LAWRENCE HIGH SCHOOL A JESUIT CHRISTIAN MINORITY INSTITUTION



<u>STUDY MATERIAL-4</u> SUBJECT – STATISTICS

Pre-test

Chapter: BIVARIATE ANALYSIS

Topic:Regression

Class: XII

Date: 13.05.2020

REGRESSION & RANK CORRELATION

PART 4

The slope of the regression line y on x is $m_1 = b_{yx}$ and the slope of the regression line x on y is $m_2 = \frac{1}{b_{yy}}$.

So the angle between the lines is $\theta = tan^{-1}\left(\frac{m_2 - m_1}{1 + m_2 m_1}\right) = tan^{-1}\left(\frac{1 - b_{yx}b_{xy}}{b_{yx} + b_{xy}}\right) =$

$$tan^{-1}(\frac{1-r_{xy}^2}{b_{yx}+b_{xy}})$$

<u>Case 1</u>: When the lines are parallel $\theta = 0 \implies 1 - r_{xy}^2 = 0 \implies r_{xy} = \pm 1$. In this case the lines coincide on each other.

<u>Case 2</u>: When the lines are perpendicular $b_{yx} + b_{xy} = 0$ $\Rightarrow b_{yx} = -b_{xy}$.

Which is not possible as both the regression coefficients should be of same sign. It is mathematically true only when $r_{xy} = 0$.

(In statistics we do not find the regression line for linearly independent variables. So we say mathematically true.)

> The mean of error is zero.

Pf: From regression line y on x, $Y_i = \overline{y} + b_{yx} (x - \overline{x})$ $\Rightarrow \sum Y_i = \sum y_i + b_{yx} \sum (x_i - \overline{x}) \Rightarrow n\overline{Y} = n \overline{y} \Rightarrow \overline{Y} = \overline{y}$

So mean of the predicted value of y is same as the corresponding observed value of y. Now $\bar{e} = \frac{1}{n} \sum e_i = \frac{1}{n} \sum (y_i - Y_i) = \bar{Y} - \bar{y} = 0$

Variance of error.

$$s_e^2 = \frac{1}{n} \sum e_i^2 - \bar{e}^2 = \frac{1}{n} \sum e_i^2 \text{ as } \bar{e} = 0$$

$$\implies s_e^2 = \frac{1}{n} (y_i - Y_i)^2 = \frac{1}{n} \sum ((y_i - \bar{y}) - r_{xy} \frac{s_y}{s_x} (x_i - \bar{x}))^2$$

$$\begin{aligned} &= \frac{1}{n} \sum (y_i - \bar{y})^2 - 2 r_{xy} \frac{s_y}{s_x} \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) + r_{xy}^2 \frac{s_y^2}{s_x^2} \frac{1}{n} \sum (x_i - \bar{x})^2 \\ &= s_y^2 - 2 r_{xy} \frac{s_y}{s_x} \cos(x, y) + r_{xy}^2 \frac{s_y^2}{s_x^2} s_x^2 \\ &= s_y^2 - 2 r_{xy} \frac{s_y}{s_x} s_x s_y r_{xy} + r_{xy}^2 s_y^2 \\ &= s_y^2 - r_{xy}^2 s_y^2 = s_y^2 (1 - r_{xy}^2) \\ &\implies s_e = s_y \sqrt{(1 - r_{xy}^2)} \end{aligned}$$

 s_e is known as standard of estimate of y in regression line y on x.

Since
$$s_e \ge 0 \implies s_y \sqrt{(1 - r_{xy}^2)} \ge 0 \implies \sqrt{(1 - r_{xy}^2)} \ge 0$$

$$\implies (1 - r_{xy}^2) \ge 0 \implies r_{xy}^2 \le 1 \implies -1 \le r_{xy} \le 1.$$

From $s_e^2 = s_y^2 (1 - r_{xy}^2)$, we see that if $r_{xy} = 0$, then $s_e^2 = s_y^2$, so that errors of estimation are as much variable as the observed values, and hence the regression line is of no help as a prediction formula. This is also seen from the fact that when $r_{xy} = 0$, we get $Y = \bar{y}$, a constant.

It is also noted that the numerical value of r_{xy} increases, s_e^2 decreases and when $r_{xy} = \pm 1$, $s_e^2 = 0$ which implies that, for each i.

$$e_i = 0 \text{ or } y_i = Y_i$$

So that all the points in scattered diagram lie on the regression line and, thus, the regression lines become a perfect prediction formula. From these observations, it is clear that the numerical value of r_{xy} can be taken as a measure of the efficiency of the regression equation as a prediction formula.

We get the same result if we consider the variance of the errors of estimates of x from its linear regression on y.

> Variance of predicted values of y in regression y on x.

$$s_Y^2 = \frac{1}{n} \sum (Y_i - \bar{Y})^2 = \frac{1}{n} \sum (\bar{y} + b_{yx}(x_i - \bar{x}) - \bar{y})^2$$
$$= b_{yx}^2 \frac{1}{n} \sum (x_i - \bar{x})^2 = r_{xy}^2 \frac{s_y^2}{s_x^2} s_x^2 = r_{xy}^2 s_y^2$$

 $\Rightarrow r_{xy}^2 = \frac{s_Y^2}{s_y^2}.$

Which is known as coefficient of determination. It may be used as measure of usefulness of linear regression equation on x.

 $|r_{xy}| = \frac{s_Y}{s_y}$. It shows the proportion of total variability of y which is accounted for by its linear regression on x.

It can similarly be shown that $|r_{xy}| = \frac{s_X}{s_x}$ where s_X is the standard deviation of the predicted values of x from its regression on y.

$$\succ Cov(x,e) = \frac{1}{n} \sum (x_i - \bar{x})e_i = \frac{1}{n} \sum x_i e_i - \bar{x} \frac{1}{n} \sum e_i = \frac{1}{n} \sum x_i e_i - \bar{x}\bar{e}_i$$
$$= \frac{1}{n} \sum x_i e_i = 0$$
$$\Rightarrow r_{xe} = 0$$
Since $\bar{e} = 0$ and $\sum x_i e_i = \sum x_i (y_i - a - bx_i)$
$$= \sum x_i y_i - na \sum x_i - b \sum x_i^2$$
$$= 0 \text{ (from 2}^{nd} \text{ normal equation)}$$

$$Cov (Y, e) = cov(a + bx, e) = cov(a, e) + b. cov(x, e) = 0$$

Since a=constant and $cov(x, e) = 0$
So $r_{Ye} = 0$

This shows that correlation between y and its predicted value Y is non-negative and numerically equal to the correlation between y and x.

Spearman's rank correlation(no tie case)

Spearman's rank correlation coefficient $r_R = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$

Where $d_i = u_i - v_i$

 u_i = Rank given by the judge 1 to the ith candidate

 v_i = Rank given by the judge 2 to the ith candidate

n = no of candidates

Example1.

The following are the results of 5 candidates in a singing competition given by two different judges.

J1:	12	18	11	8	10
J2:	15	12	9	11	13
u_i :	2	1	3	5	4
v_i :	1	3	5	4	2
d_i :	1	-2	-2	1	2
d_{i}^{2} :	0	4	4	1	4

So
$$\sum d_i^2 = 13$$
 and $r_R = 1 - \frac{6\sum d_i^2}{n(n^2-1)} = 1 - \frac{6X\,13}{5X24} = \frac{7}{20}$
Derivation of Spearman's rank correlation coefficient eith no tier that
From the above example it is clear that u_i and v_i posses the natural numbers.
So $\sum_{i=1}^n u_i^2 = \sum_{i=1}^n v_i^2 = \frac{n(n+1)}{2} \Rightarrow \bar{u} = \bar{v} = \frac{n+1}{2}$
 $\sum_{i=1}^n u_i^2 = \sum_{i=1}^n v_i^2 = \frac{n(n+1)(2n+1)}{6}$
 $s_u^2 = \frac{1}{n} \sum_{i=1}^n u_i^2 - (\bar{u})^2 = \frac{1}{n} \frac{n(n+1)(2n+1)}{6} - (\frac{n+1}{2})^2 = \frac{n^2-1}{12} = s_v^2$
 $\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (u_i - v_i)^2 = \sum_{i=1}^n u_i^2 + \sum_{i=1}^n v_i^2 - 2\sum_{i=1}^n u_i v_i$
 $\Rightarrow \frac{1}{n} \sum_{i=1}^n u_i v_i = \frac{(n+1)(2n+1)}{6} - \frac{1}{2n} \sum_{i=1}^n d_i^2$
 $\Rightarrow cov(u, v) = \frac{1}{n} \sum_{i=1}^n u_i v_i - \bar{u}\bar{v} = \frac{(n+1)(2n+1)}{6} - \frac{1}{2n} \sum_{i=1}^n d_i^2 - (\frac{n+1}{2})^2$
 $\Rightarrow cov(u, v) = \frac{(n^2-1)}{12} - \frac{1}{2n} \sum_{i=1}^n d_i^2$
So $r_R = r_{uv} = \frac{cov(u,v)}{s_u s_v} = \frac{\frac{(n^2-1)}{(\frac{n^2-1}{12})} - \frac{1}{2n} \sum_{i=1}^{n} d_i^2}{1}$

Property of rank correlation coefficient:

- ➤ The rank correlation cefficient lies between -1 and +1. $s_{d}^{2} = \frac{1}{n} \sum_{i=1}^{n} d_{i}^{2} (\bar{d})^{2} = \frac{1}{n} \sum_{i=1}^{n} d_{i}^{2} \dots (*) \text{ (since } \bar{d} = \bar{u} \bar{v} = 0)$ $s_{d}^{2} = s_{u-v}^{2}$ $\Rightarrow s_{d}^{2} \leq s_{u-v}^{2} + s_{u+v}^{2} \text{ since } s_{u+v}^{2} \geq 0$ $\Rightarrow s_{d}^{2} \leq 2(s_{u}^{2} + s_{v}^{2}) \Rightarrow s_{d}^{2} \leq 2 \frac{2(n^{2}-1)}{12} = \frac{n^{2}-1}{3} \dots (**)$ Combining (*) and (**) we get, $\frac{1}{n} \sum_{i=1}^{n} d_{i}^{2} \leq \frac{n^{2}-1}{3} \Rightarrow r_{R} \geq 1 \frac{6n(n^{2}-1)}{3n(n^{2}-1)} \Rightarrow r_{R} \geq -1 \dots (***)$ $r_{R} = 1 \frac{6 \sum d_{i}^{2}}{n(n^{2}-1)} \leq 1 \text{ since } 1 \frac{6 \sum d_{i}^{2}}{n(n^{2}-1)} \geq 0 \dots (****)$ Combining (***) and (****), $-1 \leq r_{R} \leq +1$
- ➢ In case of **perfect agreement** between two judges $u_i = v_i \Rightarrow d_i = 0$

$$\Rightarrow \sum_{i=1}^{n} d_i^2 = 0 \Rightarrow r_R = 1 - \frac{6\sum d_i^2}{n(n^2 - 1)} = 1$$

In case of **perfect disagreement** between two judges $u_i + v_i = n + 1$

$$\sum_{i=1}^{n} d_i^2 = \sum_{i=1}^{n} (u_i - v_i)^2 = \sum_{i=1}^{n} (n+1-2v_i)^2$$
$$= \sum_{i=1}^{n} (n+1)^2 + 4 \sum_{i=1}^{n} v_i^2 - 4(n+1) \sum_{i=1}^{n} v_i$$
$$= \frac{1}{3} n(n^2 - 1)$$

So
$$r_R = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} = -1$$

These two are the extreme cases . So all other cases lie in between these two case. Hence $-1 \le r_R \le +1$

- Uses of rank correlation:
 - Rank correlation is used to measure the extent of association that may exist between two series of ranks for the same set of individuals.
 - Attributes cannot be measured but individuals can be ranked according to their degree of possession of an attribute. Hence rank correlation can be used to see whether two attributes are correlated or not.
 - In case of two judges, rank correlation shows the extent of agreement between the judges.
 - Sometimes rank correlation is used instead of simple correlation because rank correlation coefficient calculation is much easier.

Prepared by Snjay Bhattacharya