# Measures of Correlation

In linear regression, there is a dependent variable Y and an independent regressor variable X. We think of Y as being approximately a function of X. We want to know if there is a relation of the form

$$E(Y|X=x) = a + bx$$

$$\underset{y}{\overset{\mu}{\sum}} y \, P(Y=y|X=x)$$

For example $\quad Y = a + bX + \varepsilon$

$$\text{where } \varepsilon \sim N(0, \sigma^2)$$

or even $\quad Y = a + bX$

In these cases $\quad Cov(X, Y) = Cov(X, a+bX)$

$$= b \, Var(X)$$

$$b = \frac{Cov(X,Y)}{Var(X)}.$$

We want estimates

of $\quad Cov(X, Y)$.

The covariance of $X$ and $Y$ is

$$Cov(X,Y) = E\left[(X - E(X))(Y - E(Y))\right]$$

$$= E[XY] - E[X]E[Y]$$

The correlation

of $X$ and $Y$ is defined by

$$Corr(X,Y) = \frac{Cov(X,Y)}{\sqrt{Var(X)Var(Y)}}$$

we have

① $-1 \le Corr(X,Y) \le +1$

② If $X$ and $Y$ are independent, $Cov(X,Y) = 0$ but $Cov(X,Y) = 0$ does not imply that $X$ and $Y$ are independent.

③ If $Y = aX + b$, then

$$Corr(X,Y) = \frac{b\,Var(X)}{\sqrt{Var(X) \cdot b^2 Var(X)}}$$

$$= \frac{b}{|b|}$$

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

We want measures of correlation
for samples $\vec{x} = (x_1, ..., x_n)$
and $\vec{y} = (y_1, ..., y_n)$.

~~The~~ Three such measures are in use.

① Pearson's Correlation Coefficient

Let $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ , $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

$$S_{xx} = \sum_{i=1}^{n} (x_i - \bar{x})^2 = \sum_{i=1}^{n} x_i^2 - \frac{\left(\sum_{i=1}^{n} x_i\right)^2}{n}$$

$$S_{yy} = \sum_{i=1}^{n} (y_i - \bar{y})^2 = \sum_{i=1}^{n} y_i^2 - \frac{\left(\sum_{i=1}^{n} y_i\right)^2}{n}$$

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

Pearson Corr. Coeff. =

$$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$$

The code for calculating $r$

$\text{Cor} (\ \langle \text{Data} ~~\rangle,~~ \text{method} = \text{"pearson"})$

or $\text{Cor} (\langle \text{Data} \rangle)$

3

## ② Spearman's Correlation coefficient

First we define the $ranks$ of the variables.

The rank of $X_i$ is $j$ if $X_i$ is the $j$th smallest of $X_1, X_2, ..., X_n$.

The rank of $Y_i$ is $j$ if $Y_i$ is the $j$th smallest of $Y_1, ..., Y_n$

Spearman's correlation coefficient $r_s$ is Pearson's correlation coefficient applied to

$$(rank(X_1), rank(X_2), ..., rank(X_n))$$
$$(rank(Y_1), rank(Y_2), ..., rank(Y_n))$$

The R code for calculating $r_s$ is

$$cor(\,<Data>,\ method = "spearman"\,)$$

# Fact

If the $x_i$'s are unique (no repetitions)
and the $y_i$'s are themselves unique
then

$$r_s = 1 - \frac{6 \sum_{i=1}^{n} d_i^2}{n(n^2-1)}$$

where $d_i = \text{rank}(y_i) - \text{rank}(x_i)$

③ Kendall's correlation coefficient

We consider ~~pairs of obser~~ vectors of observations

$$\binom{x_i}{y_i} \quad i = 1, \dots, n$$

A pair of vectors $\binom{x_i}{y_i}$ $\binom{x_j}{y_j}$

is concordant if $\quad x_i < x_j$

if either $x_i < x_j$ and $y_i < y_j$

or $x_i > x_j$ and $y_i > y_j$

A pair of vectors is discordant if

either $x_i < x_j$ and $y_i > y_j$

or $\quad x_i > x_j$ and $y_i < y_j$

Let $n_c$ be the number of concordant pairs

and $n_d$ be the number of discordant pairs.

Note: $n_c + n_d = \binom{n}{2}$

Kendall's correlation coefficient $\tau$ "tau"

is

$$\tau = \frac{n_c - n_d}{\binom{n}{2}} = \frac{n_c - n_d}{\frac{n(n-1)}{2}}$$

"Kendall's $\tau$ can be calculated from the ranks.

We can calculate it in R using

Cor ($\langle$Data$\rangle$, method = "Kendall")

# Inference on correlation coefficients

## Pearson

We need an assumption about the joint distribution of $X$ and $Y$. We assume $(X, Y)$ is bivariate normal distributed with parameters $\mu_X, \sigma_X, \mu_Y, \sigma_Y, \rho$ where $\rho$ is the correlation.

We have $\cancel{E(Y|X=x) = \mu_Y + \rho \frac{\sigma_X}{\sigma_Y}}$

$$E(Y | X=x) = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X)$$

## Theorem

Under $H_0: \rho = 0$ the statistic

$$\frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$ is $t$-distributed with $n-2$ degrees of freedom

where $r$ is Pearson's correlation $n$ is number of observations.

In R Cor.test (Data, method = "pearson").

7