

# Correlation Theory

Still bi-variate statistics

$X \sim$  random variable

$Y \sim$  random variable

(c) Pongsa Pornchaiwiseskul,

1

## Covered Topics

- Pearson's Correlation
- Spearman's Rank Correlation

(c) Pongsa Pornchaiwiseskul,

2

# Population Covariance (1)

## Definition

$$\begin{aligned}\sigma_{XY} &= E[(X - \mu_X)(Y - \mu_Y)] \\ &= \iint (x - \mu_X)(y - \mu_Y)f(x, y)dxdy\end{aligned}$$

a constant

(c) Pongsa Pornchaiwiseskul,

3

# Population Covariance (2)

## Sign of Covariance

Positive  $\implies$  if one RV is above or below its mean, the other RV tends to be also above or below its mean

Negative  $\implies$  if one RV is above or below its mean, the other RV tends to be below or above its mean

(c) Pongsa Pornchaiwiseskul,

4

# Population Covariance (3)

## Magnitude of Covariance

unbounded

depends on the units of both RV's

## Unit of covariance

= unit of X times unit of Y

e.g., X is in Baht and Y is in Kilogram

$\sigma_{XY}$  is in Baht-Kilogram

(c) Pongsa Pornchaiwiseskul,

5

# Population Correlation (1)

- Definition

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

- Sign of Correlation

—same as that of Covariance

(c) Pongsa Pornchaiwiseskul,

6

# Population Correlation (2)

## Magnitude of Correlation

always bounded between -1 and 1

$$-1 \leq \rho_{XY} \leq +1$$

## Unit of Correlation

no unit

comparable between populations

(c) Pongsa Pornchaiwiseskul,

7

# Population Correlation (3)

## Interpretation of Correlation

$\rho_{XY} = +1 \implies$  If a variable is above or below its mean, the other will be above or below its own mean with certainty

$\rho_{XY} = -1 \implies$  If a variable is above or below its mean, the other will be below or above its own mean with certainty

$\rho_{XY} = 0 \implies$  If a variable is deviated from its mean, the other will be expected at its mean

(c) Pongsa Pornchaiwiseskul,

8

# Sample Covariance

$s_{XY}$  is an estimator for  $\sigma_{XY}$

Required paired sample

Estimator

$$s_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$$

(c) Pongsa Pornchaiwiseskul,

9

## Paired Sample of Size $n$

$i$	$X_i$	$Y_i$
1	$X_1$	$Y_1$
2	$X_2$	$Y_2$
:	:	:
:	:	:
$n$	$X_n$	$Y_n$

(c) Pongsa Pornchaiwiseskul,

10

# Sample Correlation (1)

$r_{XY}$  is an estimator of  $\rho_{XY}$

Definition 
$$r_{XY} = \frac{S_{XY}}{S_X S_Y}$$

Sign of sample Correlation

same as that of sample Covariance

(c) Pongsa Pornchaiwiseskul,

11

# Sample Correlation (2)

Magnitude of Sample Correlation

same as population correlation

always bounded between -1 and 1

$$-1 \leq r_{XY} \leq +1$$

Unit of sample Correlation

no unit

comparable between data sets

(c) Pongsa Pornchaiwiseskul,

12

# Test for Zero Correlation

$$H_0 : \rho_{XY} = 0$$

$$H_1 : \rho_{XY} \neq 0$$

Theorem

$$t_{cal} = \frac{r_{XY}}{\sqrt{\frac{1-r_{XY}^2}{n-2}}} \sim t(n-2)$$

Perform a Two-sided test.

# Test for Non-zero Correlation (1)

$$H_0 : \rho_{XY} = a, \quad a \neq 0$$

$$H_1 : \rho_{XY} \neq a$$

Theorem

$$\omega = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right),$$

$$\mu_\omega = \frac{1}{2} \ln \left( \frac{1+\rho}{1-\rho} \right)$$

## Test for Non-zero Correlation (2)

$$\omega \sim N\left(\mu_{\omega}, \frac{1}{n-3}\right)$$

$$z_{cal} = \frac{\omega - \mu_{\omega}}{\sqrt{\frac{1}{n-3}}} \sim N(0,1)$$

Perform a Two-sided test.

## Rank Correlation(1)

Two judges (A and B) are to rank  $n$  different objects (contestants)

**Question:** Are the two judges correlated?

How can similarity or dissimilarity be measured?



# Rank Correlation(2)

Spearman's Rank Correlation (sample)

$$r' = 1 - \frac{6 \sum D_i^2}{n(n^2 - 1)}$$

No definition for population rank correlation

# Rank Correlation(3)

$R_{ij}$  = rank given to object  $i$  by judge  $j$

$D_i$  = rank difference for object  $i$

$$= R_{iA} - R_{iB}$$

# Rank Correlation(4)

## Paired Sample of Size $n$

$i$	$RA_i$	$RB_i$
1	$RA_1$	$RB_1$
2	$RA_2$	$RB_2$
:	:	:
:	:	:
$n$	$RA_n$	$RB_n$

(c) Pongsa Pornchaiwiseskul,

19

# Rank Correlation(5)

## Magnitude of Correlation

always bounded between -1 and 1

$$-1 \leq r'_{XY} \leq +1$$

## Unit of Correlation

no unit

comparable between populations

(c) Pongsa Pornchaiwiseskul,

20

# Rank Correlation(6)

## Interpretation of Sample Rank Correlation

$r'_{XY} = +1 \implies$  If both judges totally agree on the rankings of all the  $n$  objects

$r'_{XY} = -1 \implies$  If both judges totally disagree on the rankings of all the  $n$  objects

$r'_{XY} = 0 \implies$  If the two judges are uncorrelated

## Test for Zero Rank Correlation

$$H_0 : \rho'_{XY} = 0$$

$$H_1 : \rho'_{XY} \neq 0$$

Theorem

$$z_{cal} = \frac{r'_{XY}}{\sqrt{\frac{1}{n-1}}} \sim N(0,1)$$

Perform a Two-sided Z-test.

# Test for Non-zero Rank Correlation

No such a thing??

(c) Pongsa Pornchaiwiseskul,

23

## Correlation Theory

Now tri-variate

$X \sim$  random variable

$Y \sim$  random variable

$Z \sim$  random variable

(c) Pongsa Pornchaiwiseskul,

1

# Covered Topics

- Partial Correlation (Pearson)

## Partial Correlation (1)

- X, Y and Z are three RV's
- They are assumed to be related.
- Ignoring Z,  $\text{Corr}(X, Y) = \text{“direct”}$   
correlation (X, Y) + indirect effects from Z

# Partial Correlation (2)

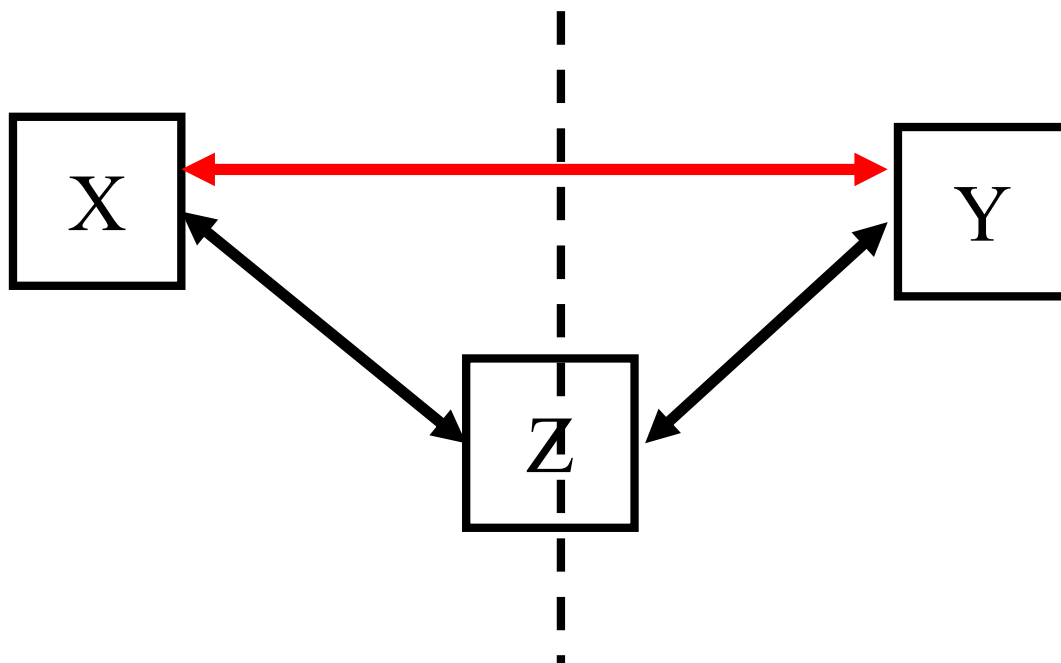
## Partial Correlation

- Correlation when embedded effect of  $Z$  has been removed from both  $X$  and  $Y$
- Partial Corr ( $X, Y$ ) = “direct” corr( $X, Y$ )

(c) Pongsa Pornchaiwiseskul,

4

# Partial Correlation (3)



(c) Pongsa Pornchaiwiseskul,

5

# Partial Correlation (4)

Assumptions

- X.Z are related as  $X = \alpha_1 + \alpha_2 Z + \varepsilon$
- Y.Z are related as  $Y = \beta_1 + \beta_2 Z + \xi$

# Partial Correlation (5)

PCorr

- X.Z are related as  $X = \alpha_1 + \alpha_2 Z + \varepsilon$
- Y.Z are related as  $Y = \beta_1 + \beta_2 Z + \xi$

# Partial Correlation (6)

By OLS

$$\hat{\alpha}_2 = \frac{\sum x_i z_i}{\sum z_i^2}, \hat{\alpha}_1 = \bar{X} - \hat{\alpha}_2 \bar{Z}$$

where  $x_i = X_i - \bar{X},$

$$z_i = Z_i - \bar{Z}$$

(c) Pongsa Pornchaiwiseskul,

8

# Partial Correlation (7)

By OLS

$$\hat{\beta}_2 = \frac{\sum y_i z_i}{\sum z_i^2}, \hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{Z}$$

where  $y_i = y_i - \bar{Y},$

$$z_i = Z_i - \bar{Z}$$

(c) Pongsa Pornchaiwiseskul,

9



## Partial Correlation (8)

$$X \text{ without } Z \quad X - \alpha_2 Z = \alpha_1 + \varepsilon$$

$$Y \text{ without } Z \quad Y - \beta_2 Z = \beta_1 + \xi$$

Partial Corr between X and Y ( $\rho_{XY.Z}$ )

$$= \text{corr}(\alpha_2 + \varepsilon, \beta_1 + \xi) = \text{corr}(\varepsilon, \xi)$$

## Partial Correlation (9)

Sample partial correlation coefficient

$$r_{XY.Z} = \frac{\sum (X_i - \hat{X}_i)(Y_i - \hat{Y}_i)}{\sqrt{\sum (X_i - \hat{X}_i)^2} \sqrt{\sum (Y_i - \hat{Y}_i)^2}}$$

# Partial Correlation (9)

Sample partial correlation coefficient

$$r_{XY.Z} = \frac{\sum (x_i - \hat{\alpha}_2 z_i)(y_i - \hat{\beta}_2 z_i)}{\sqrt{\sum (x_i - \hat{\alpha}_2 z_i)^2} \sqrt{\sum (y_i - \hat{\beta}_2 z_i)^2}}$$

# Partial Correlation (9)

Sample partial correlation coefficient

$$\begin{aligned} \sum (x_i - \hat{\alpha}_2 z_i)^2 &= \sum x_i^2 - 2\hat{\alpha}_2 \sum x_i z_i + \hat{\alpha}_2^2 \sum z_i^2 \\ &= \sum x_i^2 - 2 \frac{\sum x_i z_i}{\sum z_i^2} \sum x_i z_i + \left( \frac{\sum x_i z_i}{\sum z_i^2} \right)^2 \sum z_i^2 \\ &= \sum x_i^2 - \frac{(\sum x_i z_i)^2}{\sum z_i^2} \end{aligned}$$

# Partial Correlation (10)

Sample partial correlation coefficient

$$\begin{aligned} &= \sum x_i^2 - \frac{\left(\sum x_i z_i\right)^2}{\sum x_i^2 \sum z_i^2} \sum x_i^2 \\ &= \left(1 - r_{XZ}^2\right) \sum x_i^2 \end{aligned}$$

# Partial Correlation (11)

Sample partial correlation coefficient

$$\sum \left(y_i - \hat{\alpha}_2 z_i\right)^2 = \left(1 - r_{YZ}^2\right) \sum y_i^2$$

# Partial Correlation (12)

Sample partial correlation coefficient

$$\begin{aligned} & \sum (x_i - \hat{\alpha}_2 z_i)(y_i - \hat{\beta}_2 z_i) \\ &= \sum x_i y_i - \hat{\alpha}_2 \sum y_i z_i \\ & \quad - \hat{\beta}_2 \sum x_i z_i + \hat{\alpha}_2 \hat{\beta}_2 \sum z_i^2 \end{aligned}$$

(c) Pongsa Pornchaiwiseskul,

16

# Partial Correlation (13)

Sample partial correlation coefficient

$$\begin{aligned} &= \sum x_i y_i - \frac{\sum x_i z_i}{\sum z_i^2} \sum y_i z_i \\ & \quad - \frac{\sum y_i z_i}{\sum z_i^2} \sum x_i z_i + \frac{\sum x_i z_i}{\sum z_i^2} \frac{\sum y_i z_i}{\sum z_i^2} \sum z_i^2 \\ &= \sum x_i y_i - \frac{\sum y_i z_i}{\sum z_i^2} \sum x_i z_i \end{aligned}$$

(c) Pongsa Pornchaiwiseskul,

17

# Partial Correlation (14)

Sample partial correlation coefficient

$$\begin{aligned} &= r_{YX} \sqrt{\sum x_i^2} \sqrt{\sum y_i^2} \\ &\quad - \frac{r_{YZ} \sqrt{\sum y_i^2} \sqrt{\sum z_i^2}}{\sum z_i^2} r_{XZ} \sqrt{\sum x_i^2} \sqrt{\sum z_i^2} \\ &= (r_{YX} - r_{YZ} r_{XZ}) \sqrt{\sum x_i^2} \sqrt{\sum y_i^2} \end{aligned}$$

(c) Pongsa Pornchaiwiseskul,

18

# Partial Correlation (15)

Sample partial correlation coefficient

$$\begin{aligned} r_{XY.Z} &= \frac{(r_{YX} - r_{YZ} r_{XZ}) \sqrt{\sum x_i^2} \sqrt{\sum y_i^2}}{\sqrt{(1 - r_{XZ}^2) \sum x_i^2} \sqrt{(1 - r_{YZ}^2) \sum y_i^2}} \\ &= \frac{r_{YX} - r_{YZ} r_{XZ}}{\sqrt{1 - r_{XZ}^2} \sqrt{1 - r_{YZ}^2}} \end{aligned}$$

(c) Pongsa Pornchaiwiseskul,

19