

## Multi-variate Statistics

Extension of Bi-variate Statistics

$(Y, \mathbf{X}) \sim$  random variables

where

$\mathbf{X} \sim$  vectors of  $K$  random variables

$$\mathbf{X} = [X_1, X_2, \dots, X_K]$$

$Y \sim$  a single random variable

## Multiple Regression Analysis

Focus on the dependency of  $Y$  on the  $\mathbf{X}$  vector, e.g.,

$$\mu_{Y|\mathbf{X}} = m(X_1, X_2, \dots, X_K) = m(\mathbf{X})$$

$$\sigma_{Y|\mathbf{X}}^2 = v(X_1, X_2, \dots, X_K) = v(\mathbf{X})$$

$X_k$  - explanatory or independent variable,

$$k = 1, \dots, K$$

$Y$  - dependent variable

## Multi-variate Analyses

- Pair-wise Covariance or Correlation
- Multi-way ANOVA
- Multiple Regression

## Multiple Linear Regression

### Assumptions

1) linearity  $\mu_{Y|\mathbf{X}} = \mathbf{X}\boldsymbol{\beta}$

where  $\boldsymbol{\beta} = [\beta_1 \ \beta_2 \ \dots \ \beta_K]^T$  are unknown parameters

2) variance-independent or  $\sigma_{Y|\mathbf{X}}^2 = \sigma^2$

3) normality, i.e.  $Y|\mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2)$

## CLNRM (1)

### Classical Linear Normal Regression

Model is based upon the assumptions

$$Y_i = X_i \beta + \varepsilon_i$$

where  $i$  = index of the observation

$\varepsilon_i$  = identical and independent

normal error term

$\varepsilon_i \sim N(0, \sigma^2)$  for all  $i=1, \dots, n$

## CLNRM

### Matrix Representation (1)

Define

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} X_{11} & X_{21} & \dots & X_{K1} \\ X_{12} & X_{22} & \dots & X_{K1} \\ \vdots & \vdots & \vdots & \vdots \\ X_{1n} & X_{2n} & \dots & X_{Kn} \end{bmatrix}, \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

## CLNRM (2)

$X_i$  is pre-selected or non-random but  $Y_i$  or

$\varepsilon_i$  is randomly sampled.

$X_i \beta$  is the non-random component of  $Y_i$

$\varepsilon_i$  is the random component of  $Y_i$ .

Note that  $X_1$  can be intentionally set to one for all observations so that its coefficient  $\beta_1$  becomes the y-intercept.

## CLNRM

### Matrix Representation (2)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim MVN(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

where

$\mathbf{0}$  is a  $n \times 1$  column vector of zeroes

$\mathbf{I}_n$  is an  $n \times n$  identity matrix.

## CLNRM

### Matrix Representation (3)

$\mathbf{X}$  is non-random. It is required that the matrix  $\mathbf{X}^T\mathbf{X}$  is invertible. Why?

Remember why we need  $\sum_{i=1}^n (X_i - \bar{X})^2 > 0$  in Simple Linear Regression?

## OLS Estimation for CLNRM (2)

First-Order Conditions

$$2[-\mathbf{X}]^T [\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}] = \mathbf{0}$$

$$-\mathbf{X}^T\mathbf{Y} + \mathbf{X}^T\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$$\hat{\boldsymbol{\beta}} = [\mathbf{X}^T\mathbf{X}]^{-1} \mathbf{X}^T\mathbf{Y}$$

## OLS Estimation for CLNRM (1)

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n [Y_i - (X_{1i}\beta_1 + X_{2i}\beta_2 + \dots + X_{Ki}\beta_K)]^2$$

or

$$\min_{\boldsymbol{\beta}} [\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}]^T [\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}]$$

## OLS Estimation for CLNRM (3)

Estimator for  $\sigma^2$

$$\hat{\sigma}^2 = \frac{1}{n-K} [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}]^T [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}]$$

$$= \frac{1}{n-K} [\mathbf{Y}^T\mathbf{Y} - \mathbf{Y}^T\hat{\mathbf{Y}}]$$

where  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  is called the fitted value of  $\mathbf{Y}$

Why  $n-K$ ?

## Properties of OLS estimators (1)

Theorem  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$

$$V(\hat{\boldsymbol{\beta}}) = \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}$$

Does not require normality assumption.

Note that  $\hat{\boldsymbol{\beta}}$  is an unbiased estimator of  $\boldsymbol{\beta}$ .

## Properties of OLS estimators (3)

Proof  $V(\hat{\boldsymbol{\beta}}) = [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T V(\mathbf{Y}) [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T V(\mathbf{Y}) \mathbf{X} [\mathbf{X}^T \mathbf{X}]^{-1}$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T V(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \mathbf{X} [\mathbf{X}^T \mathbf{X}]^{-1}$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T V(\boldsymbol{\varepsilon}) \mathbf{X} [\mathbf{X}^T \mathbf{X}]^{-1}$   
 $= \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{I}_n \mathbf{X} [\mathbf{X}^T \mathbf{X}]^{-1}$   
 $= \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}$

## Properties of OLS estimators (2)

Proof  $E(\hat{\boldsymbol{\beta}}) = [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T E(\mathbf{Y})$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T E(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T [\mathbf{X}\boldsymbol{\beta} + E(\boldsymbol{\varepsilon})]$   
 $= [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}$   
 $= \boldsymbol{\beta}$

## Properties of OLS estimators (4)

Theorem Due to the normality assumption

of  $\boldsymbol{\varepsilon}$ ,

$$\hat{\boldsymbol{\beta}} \sim MVN\left(\boldsymbol{\beta}, \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}\right)$$

and  $(n - K) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n - K)$

## Properties of OLS estimators (5)

Variance-Covariance Matrix of  $\hat{\beta}$

$$V(\hat{\beta}) = \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}$$

$$= \begin{bmatrix} V(\hat{\beta}_1) & C(\hat{\beta}_1, \hat{\beta}_2) & \cdots & C(\hat{\beta}_1, \hat{\beta}_K) \\ C(\hat{\beta}_2, \hat{\beta}_1) & V(\hat{\beta}_2) & \cdots & C(\hat{\beta}_2, \hat{\beta}_K) \\ \vdots & \vdots & \ddots & \vdots \\ C(\hat{\beta}_K, \hat{\beta}_1) & C(\hat{\beta}_K, \hat{\beta}_2) & \cdots & V(\hat{\beta}_K) \end{bmatrix}$$

$\sigma^2$  is generally unknown.

## Properties of OLS estimators (7)

Standard Deviation of  $\hat{\beta}_k$

$$sd(\hat{\beta}_k) = \sqrt{V(\hat{\beta}_k)}$$

Standard Error of  $\hat{\beta}_k$

$$se(\hat{\beta}_k) = \sqrt{\hat{V}(\hat{\beta}_k)}$$

## Properties of OLS estimators (6)

Estimated Variance-Covariance Matrix of  $\hat{\beta}$

$$\hat{V}(\hat{\beta}) = \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

$$= \begin{bmatrix} \hat{V}(\hat{\beta}_1) & \hat{C}(\hat{\beta}_1, \hat{\beta}_2) & \cdots & \hat{C}(\hat{\beta}_1, \hat{\beta}_K) \\ \hat{C}(\hat{\beta}_2, \hat{\beta}_1) & \hat{V}(\hat{\beta}_2) & \cdots & \hat{C}(\hat{\beta}_2, \hat{\beta}_K) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{C}(\hat{\beta}_K, \hat{\beta}_1) & \hat{C}(\hat{\beta}_K, \hat{\beta}_2) & \cdots & \hat{V}(\hat{\beta}_K) \end{bmatrix}$$

## Properties of OLS estimators (8)

$$t_{cal} = \frac{\hat{\beta}_k - \beta_k}{sd(\hat{\beta}_k)} = \frac{\hat{\beta}_k - \beta_k}{\sqrt{\frac{\hat{\sigma}^2}{(n-K)} \frac{\sigma^2}{\sigma^2}}} = \frac{\hat{\beta}_k - \beta_k}{se(\hat{\beta}_k)} \sim t(n-K)$$

<<Basis for statistical inference>>

## Central Limit Theorem (1)

Similar to that for the Simple Linear Regression Model. Even though the error terms are not normal, the properties of OLS estimators asymptotically hold when the sample size is very large.

## Gauss-Markov Theorem (1)

Similar to that for the Simple Linear Regression Model. Given that  $\mathbf{X}$  is non-random, OLS estimator is Best Linear Unbiased Estimator.

## Central Limit Theorem (2)

In mathematical term,

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \overset{A}{\sim} \text{MVN}\left(\mathbf{0}, \sigma^2 \left[ \frac{\mathbf{X}^T \mathbf{X}}{n} \right]^{-1}\right)$$

## Gauss-Markov Theorem (2)

$\hat{\boldsymbol{\beta}}$  is OLS estimator of  $\boldsymbol{\beta}$

$\tilde{\boldsymbol{\beta}}$  is a non-OLS linear unbiased estimator of  $\boldsymbol{\beta}$

$$\mathbf{hV}(\hat{\boldsymbol{\beta}})\mathbf{h}^T \leq \mathbf{hV}(\tilde{\boldsymbol{\beta}})\mathbf{h}^T$$

for any vector  $\mathbf{h} \neq \mathbf{0}$

## Coefficient of Determination (1)

$R^2$  is a measure for goodness-of-fit. How well does the model fit the observed data? Low  $R^2$  implies “bad” fit.

Definition 
$$R^2 \equiv 1 - \frac{SSR}{SST}$$

SSR = Sum of Squared Residuals

SST = Sum of Squared Totals

## Coefficient of Determination (3)

Low  $R^2$  or a bad fit does not mean a bad model. It simply implies a large uncertainty in the nature. It is mainly used as a criterion to select various “candidate” models.

## Coefficient of Determination (2)

where 
$$SSR = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = [\mathbf{Y} - \hat{\mathbf{Y}}]^T [\mathbf{Y} - \hat{\mathbf{Y}}]$$
$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Note that, in general,  $R^2$  cannot be greater than one but could be negative.

## Coefficient of Determination (4)

If an  $X_i$  has constant value or a linear combination of  $X_i$  's is equivalent to a constant value, then,  $0 \leq R^2 \leq 1$  always

and 
$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

## Coefficient of Determination (5)

Interpretation if  $0 \leq R^2 \leq 1$

$1-R^2$  or  $SSR/SST$  can be interpreted as the fraction of total variation of  $Y$  due to the random component ( $\mathcal{E}$ ).

$R^2$  is generally regarded as the fraction of total variation of  $Y$  explained by the explanatory variables or due to the non-random component.

## Adjusted- $R^2$ (2)

Definition

$$\bar{R}^2 \equiv 1 - \frac{SSR/(n-K)}{SST/(n-1)} = 1 - \frac{\widehat{\sigma}^2}{s_Y^2}$$

Concept

Penalize  $R^2$  by dividing with  $(n-K)$  when an irrelevant variable is added.

## Adjusted- $R^2$ (1)

We can cheat on  $R^2$  by adding more irrelevant independent variables on the right-hand side, especially when sample is small.

Higher  $K \implies$  smaller  $SSR \implies$  higher  $R^2$

## Adjusted- $R^2$ (3)

Purpose

For a small sample, it is a better measure for goodness-of-fit than  $R^2$ . It is also used as criterion to add or remove an explanatory variable from the model if it does not contradict theories.