

# Multiple Least Squares Regression

Christopher Ting

July 14, 2018

Course Web Site: <http://cting.x10host.com/CUHKSZ/CUHKSZ.html>

# Broad Lesson Plan

- 1 Introduction
- 2 OLS in Matrix Form
- 3 Properties of OLS estimator
- 4 OLS Inference & Forecast

# Topics Covered

- ➔ Fama-French 3-factor model
- ➔ Classical assumptions of linear regression
- ➔ Ordinary least squares in matrix form
- ➔ Projection matrix and residual matrix
- ➔ Small sample properties of OLS estimator—unbiasedness, normality, efficiency, unbiased variance of residuals
- ➔ Statistical inference and confidence region
- ➔ Analysis of variance

# Linear Models

- ➔ There are  $p$  explanatory variables  $X$  to explain the independent variable  $Y$ .
- ➔ Given  $n$  sets of observations of  $X$  and  $Y$ , the linear specification of their relationship is given by

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \cdots + \beta_p X_{i,p} + \epsilon_i,$$

where  $\epsilon_i$  is the noise, disturbances, or errors.

- ➔ There are various ways to estimate the  $y$  intercept and the  $p$  parameters  $\beta_j, j = 1, 2, \dots, p$ .
- ➔ The parameter  $\beta_j$  is also called the loading on  $X_j$ .

## Example: Fama-French 3-Factor Model

$$r_{st} - r_{ft} = \beta_0 + \beta_1 \text{MKT}_t + \beta_2 \text{SMB}_t + \beta_3 \text{HML}_t + \epsilon_t$$

### \* Market Factor

$$\text{MKT}_t = r_{mt} - r_{ft}.$$

### \* Size Factor

$$\begin{aligned} \text{SMB}_t = & \frac{1}{3} (\text{Small Value} + \text{Small Neutral} + \text{Small Growth}) \\ & - \frac{1}{3} (\text{Big Value} + \text{Big Neutral} + \text{Big Growth}). \end{aligned}$$

### \* Value Factor

$$\begin{aligned} \text{HML}_t = & \frac{1}{2} (\text{Small Value} + \text{Big Value}) \\ & - \frac{1}{2} (\text{Small Growth} + \text{Big Growth}). \end{aligned}$$

# Classical Assumptions

- ▶ Linearity of the conditional expectation

$$\mathbb{E}(Y_i | X_{i,1}, X_{i,2}, \dots, X_{i,p}) = \beta_0 + \beta_1 X_{i,1} + \dots + \beta_p X_{i,p}.$$

- ▶ Independent noise

$$\epsilon_1, \epsilon_2, \dots, \epsilon_n \text{ are i.i.d.}$$

with zero mean,  $\mathbb{E}(\epsilon_i)$  for all  $i$ .

- ▶ Constant variance

$$\mathbb{V}(\epsilon_i) = \sigma_\epsilon^2, \quad \text{for all } i.$$

- ▶ Gaussian noise

$$\epsilon_i \sim N(0, \sigma_\epsilon^2), \quad \text{for all } i.$$

# Implications of Assumptions

- ▶ The first assumption is implied by

$$\mathbb{E}(\epsilon_i | X_{i,1}, X_{i,2}, \dots, X_{i,p}) = 0.$$

- ▶ From the first assumption, we also have

$$\beta_j = \frac{\partial \mathbb{E}(Y_i | X_{i,1}, X_{i,2}, \dots, X_{i,p})}{\partial X_{i,j}}.$$

- ▶ Therefore,  $\beta_j$  is the change in the expected value of  $Y_i$  when  $X_{i,j}$  changes one unit.

# Conditional Variance of Independent Variable

- ▀ All  $n$  observations of  $p$   $X$ s are organized into the data matrix

$$\mathbf{X} = \begin{pmatrix} 1 & X_{1,1} & X_{1,2} & \dots & X_{1,p-1} & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & \dots & X_{2,p-1} & X_{2,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & X_{n,1} & \vdots & \dots & X_{n,p-1} & X_{n,p} \end{pmatrix}$$

- ▀ What is the conditional variance of  $Y_i$  given all the observations  $\mathbf{X}$ ?
- ▀ Answer: It is simply the variance of the error.

$$\mathbb{V}(Y_i | \mathbf{X}) = \sigma_\epsilon^2.$$

## In Matrix Form

- \* The independent variable is a vector of  $n$  rows, i.e.,

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}$$

- \* The parameters  $\beta_0, \beta_1, \dots, \beta_p$  is also assembled as a vector  $\beta$  of  $1 + p$  rows.

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

- \* Then

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon.$$

# Error

- \* The  $i$ -th row of the  $\mathbf{X}$  matrix corresponding to  $Y_i$  is transposed to form a column vector

$$\mathbf{x}_i = \begin{pmatrix} 1 \\ X_{i,1} \\ X_{i,2} \\ \vdots \\ X_{i,p} \end{pmatrix}$$

- \* Hence,  $\mathbf{x}_i^\top \boldsymbol{\beta}$  is a scalar equal to  $\beta_0 + \beta_1 X_{i,1} + \cdots + \beta_p X_{i,p}$
- \* Accordingly the error term at  $i$  is given by

$$\epsilon_i = Y_i - \mathbf{x}_i^\top \boldsymbol{\beta}.$$

# Ordinary Least Squares

- ❄ The ordinary least squares (OLS) estimator is defined as the value that minimizes the sum of the squared errors

$$\hat{\boldsymbol{\beta}} = \min_{\boldsymbol{\beta}} s(\boldsymbol{\beta}),$$

where

$$\begin{aligned} s(\boldsymbol{\beta}) &= \sum_{i=1}^n \epsilon_i^2 \\ &= \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{Y}^\top \mathbf{Y} - 2\mathbf{Y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} \\ &= \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \end{aligned}$$

## Vector of Parameter Estimates

- ❄ The OLS minimizes the Euclidean distance between  $\mathbf{Y}$  and  $\mathbf{X}\boldsymbol{\beta}$ .
- ❄ Differentiate  $s(\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$  yields the first-order condition

$$-2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} = \mathbf{0}.$$

- ❄ The solution is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

- ❄ To verify that this is a minimum, the sufficient condition is that the second derivative, namely,  $\mathbf{X}^\top \mathbf{X}$  is positive. Indeed, the rank of  $\mathbf{X}^\top \mathbf{X}$  is  $k := 1 + p$  and hence  $\mathbf{X}^\top \mathbf{X}$  is a positive definite  $k \times k$  matrix.

# Insight!

\* The fitted values are the vector  $\hat{Y} = X\hat{\beta}$ .

\* The residuals are the vector  $\hat{\epsilon} = Y - X\hat{\beta}$ .

\* Note that

$$Y = X\beta + \epsilon = X\hat{\beta} + \hat{\epsilon}.$$

\* Also, the first order conditions can be written as

$$X^T Y - X^T X \hat{\beta} = 0$$

$$X^T (Y - X \hat{\beta}) = 0$$

$$X^T \hat{\epsilon} = 0$$

\* The OLS residuals  $\hat{\epsilon}$  are orthogonal to  $X^T$

# Projection Matrix

- ❄  $\mathbf{X}\hat{\boldsymbol{\beta}}$  is the projection of  $\mathbf{Y}$  onto the span of  $\mathbf{X}$ , which is

$$\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- ❄ Therefore, the matrix that projects  $\mathbf{Y}$  onto the span of  $\mathbf{X}$  is

$$\mathbf{P}_X := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

# Orthogonal Projection Matrix

\* We have

$$\begin{aligned}\epsilon &= Y - X\hat{\beta} \\ &= Y - X(X^\top X)^{-1}X^\top Y \\ &= (I_n - X(X^\top X)^{-1}X^\top)Y\end{aligned}$$

\* So the matrix that projects  $Y$  onto the space orthogonal to the span of  $X$  is

$$M_X = I_n - X(X^\top X)^{-1}X^\top = I_n - P_X$$

# Insight!

✱ Since  $\hat{\boldsymbol{\epsilon}} = \mathbf{M}_X \mathbf{Y}$ , we have

$$\mathbf{Y} = \mathbf{P}_X \mathbf{Y} + \mathbf{M}_X \mathbf{Y} .$$

- ✱ These two projection matrices decompose the  $n$  dimensional vector  $\mathbf{Y}$  into two orthogonal components: the portion that lies in the  $k$ -dimensional space defined by  $\mathbf{X}$ ; and the portion that lies in the orthogonal  $(n - k)$ -dimensional space.
- ✱ Note that both  $\mathbf{P}_X$  and  $\mathbf{M}_X$  are symmetric and idempotent.

## Small Sample Properties: Unbiasedness

\* For  $\hat{\beta}$  we have

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \epsilon) \\ &= \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon,\end{aligned}$$

then applying the strong exogeneity assumption and the law of iterated expectation, we have

$$\begin{aligned}\mathbb{E}\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon\right) &= \mathbb{E}\left(\mathbb{E}\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \mid \mathbf{X}\right)\right) \\ &= \mathbb{E}\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\epsilon \mid \mathbf{X})\right) \\ &= 0.\end{aligned}$$

\* So the OLS estimator is unbiased.

## Small Sample Properties: Conditional Normality

- \*  $\hat{\beta}|\mathbf{X} \sim N(\beta, (\mathbf{X}^\top \mathbf{X})^{-1}\sigma_\epsilon^2)$ , where  $\sigma_\epsilon^2 = \mathbb{E}(\epsilon_i^2)$ .
- \* Proof: Note that  $\hat{\beta} = \beta + (\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top \epsilon$ , then  $\hat{\beta}$  is a linear function of  $\epsilon$ .
- \* It is sufficient to obtain the conditional variance.

$$\begin{aligned}
 \mathbb{V}(\hat{\beta}|\mathbf{X}) &= \mathbb{E}\left((\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top | \mathbf{X}\right) \\
 &= \mathbb{E}\left((\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top \epsilon \epsilon^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} | \mathbf{X}\right) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top \mathbb{E}(\epsilon \epsilon^\top | \mathbf{X}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1}\sigma_\epsilon^2.
 \end{aligned}$$

- \* Notice that the above results is true irrespective of the distribution of  $\epsilon$ .

## Small Sample Properties: Efficiency

- (Gauss-Markov theorem): In the classical linear regression model, the OLS estimator,  $\hat{\beta}$ , is the minimum variance linear unbiased estimator of  $\beta$ .
- Proof: Consider a new linear estimator  $\tilde{\beta} = \mathbf{C}\mathbf{Y}$ , where  $\mathbf{C}$  is given by

$$\mathbf{C} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top + \mathbf{D},$$

and  $\mathbf{D}$  is a nonzero matrix.

- If this estimator is unbiased, then we must have  $\mathbf{C}\mathbf{X} = \mathbf{I}$  since

$$\begin{aligned} \mathbb{E}(\mathbf{C}\mathbf{Y}) &= \mathbb{E}(\mathbf{C}\mathbf{X}\beta + \mathbf{C}\epsilon) \\ &= \mathbf{C}\mathbf{X}\beta \end{aligned}$$

- Therefore,  $\mathbf{C}\mathbf{X} = \mathbf{I}$ .

## Small Sample Properties: Efficiency (cont)

- The variance of an unbiased  $\tilde{\beta}$  is

$$\mathbb{V}(\tilde{\beta}) = \mathbf{C}\mathbf{C}^\top \sigma_\epsilon^2.$$

- Since  $\mathbf{D} = \mathbf{C} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  and  $\mathbf{C}\mathbf{X} = \mathbf{I}$ , then  $\mathbf{D}\mathbf{X} = \mathbf{0}$ .
- Consequently,

$$\begin{aligned} \mathbb{V}(\tilde{\beta}|\mathbf{X}) &= \left( \mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \left( \mathbf{D} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right)^\top \sigma_\epsilon^2 \\ &= \left( \mathbf{D}\mathbf{D}^\top + (\mathbf{X}^\top \mathbf{X})^{-1} \right) \sigma_\epsilon^2 \\ &= \mathbf{D}\mathbf{D}^\top \sigma_\epsilon^2 + \mathbb{V}(\hat{\beta}|\mathbf{X}) \end{aligned}$$

- Since  $\mathbf{D}\mathbf{D}^\top$  is positive definite, it follows that

$$\mathbb{V}(\tilde{\beta}) \geq \mathbb{V}(\hat{\beta}).$$



## Small Sample Properties: Unbiased $\hat{\sigma}_\epsilon^2$

- The estimate for  $\sigma_\epsilon^2$  is

$$\hat{\sigma}_\epsilon^2 = \frac{1}{n-k} \hat{\boldsymbol{\epsilon}}^\top \hat{\boldsymbol{\epsilon}} = \frac{1}{n-k} \boldsymbol{\epsilon}^\top \mathbf{M}_X \boldsymbol{\epsilon},$$

where  $\mathbf{M}_X$  is idempotent matrix that projects onto the space orthogonal to the span of  $\mathbf{X}$  and  $\hat{\boldsymbol{\epsilon}} = \mathbf{M}_X \mathbf{Y}$ .

- Proof:

$$\begin{aligned} \mathbb{E}(\hat{\sigma}_\epsilon^2) &= \frac{1}{n-k} \mathbb{E}(\text{Trace}(\boldsymbol{\epsilon}^\top \mathbf{M}_X \boldsymbol{\epsilon})) \\ &= \frac{1}{n-k} \mathbb{E}(\text{Trace}(\mathbf{M}_X \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top)) \\ &= \frac{1}{n-k} \text{Trace}(\mathbb{E}(\mathbf{M}_X \boldsymbol{\epsilon} \boldsymbol{\epsilon}^\top)) \\ &= \frac{1}{n-k} \sigma_\epsilon^2 \text{Trace}(\mathbf{M}_X) \\ &= \sigma_\epsilon^2. \end{aligned}$$

## Summary

- ◆ Parameter estimates for the linear model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  is a  $k$ -vector:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- ◆ The variance-covariance of the estimate  $\hat{\boldsymbol{\beta}}$  is

$$\mathbb{V}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma_\epsilon^2$$

## Statistical Inference

- The marginal distribution of any regression coefficient  $\hat{\beta}_j$  is normal with mean  $\beta_j$  and variance  $\sigma_\epsilon^2 c_{jj}$ , where  $c_{jj}$  is the  $j$ -th diagonal element of the matrix  $(\mathbf{X}^\top \mathbf{X})^{-1}$ .
- Since  $\hat{\beta} | \mathbf{X} \sim N(\boldsymbol{\beta}, (\mathbf{X}^\top \mathbf{X})^{-1} \sigma_\epsilon^2)$ , we have, for all  $j = 0, 1, \dots, p$ ,

$$t \text{ statistic} = \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_\epsilon \sqrt{c_{jj}}} \sim t_{n-k}$$

- The  $\alpha\%$  significance level for  $\beta_j$  is

$$\hat{\beta}_j - q \hat{\sigma}_\epsilon \sqrt{c_{jj}} \leq \beta_j \leq \hat{\beta}_j + q \hat{\sigma}_\epsilon \sqrt{c_{jj}}.$$

where  $q$  is the  $(1 - \alpha/2)$ -th quantile of the  $t_{n-k}$  distribution.

- The estimate  $\hat{\beta}_j$  is said to be statistically significant at the  $\alpha\%$  significance level if the absolute value of the computed  $t$  statistic is greater than  $q$ .

## Confidence Interval for Mean Response

- What is the confidence interval for the mean response  $\beta^\top \mathbf{x}$  for a given observation  $\mathbf{x}$ ?

- Given the unbiased estimate  $\hat{\beta}$ , the variance of  $\hat{\beta}^\top \mathbf{x}$  is

$$\mathbb{V}(\hat{\beta}^\top \mathbf{x}) = \mathbb{V}(\mathbf{x}^\top \hat{\beta}) = \mathbf{x}^\top \mathbb{C}(\hat{\beta}) \mathbf{x} = \sigma_\epsilon^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}$$

- Hence a  $100 \times (1 - \alpha)\%$  confidence interval for the mean response  $\beta^\top \mathbf{x}$  is

$$\hat{\beta}^\top \mathbf{x} \pm q\sigma_\epsilon \sqrt{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}}$$

## Prediction Interval for a New Observation

- ✿ Suppose a future observation of  $\mathbf{x}$  is obtained. Then, since  $\epsilon \sim N(0, \sigma_\epsilon^2)$ , we have

$$\mathbb{V}(y - \hat{\boldsymbol{\beta}}^\top \mathbf{x}) = \mathbb{V}(\epsilon) + \mathbb{V}(\hat{\boldsymbol{\beta}}^\top \mathbf{x})$$

- ✿ Hence a  $100 \times (1 - \alpha)\%$  prediction interval for  $y$  is

$$\hat{\boldsymbol{\beta}}^\top \mathbf{x} \pm q\sigma_\epsilon \sqrt{1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}}$$

# Confidence Region for All Regression Coefficients

- ▲ The statistic

$$\frac{(\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta) / k}{\sigma_\epsilon^2} \sim F_{k, n-k}$$

- ▲ Let  $F_\alpha$  be the  $(1 - \alpha)$ -th quantile of the  $F_{k, n-k}$  distribution. Then the region that simultaneously contains all regression coefficients with probability  $1 - \alpha$  is

$$\left\{ \beta : \frac{(\hat{\beta} - \beta)^\top (\mathbf{X}^\top \mathbf{X}) (\hat{\beta} - \beta) / k}{\sigma_\epsilon^2} \leq F_\alpha \right\}.$$

# Analysis of Variance

- The total variation in  $Y$  is the variation of  $Y$  with respect to  $\bar{Y}$  in terms of the sum of squares.

$$\text{TSS} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- It is partitioned into
  - the variation that can be explained by  $X_1, \dots, X_p$ . The explained sum of squares is

$$\text{ESS} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

- the variation that cannot be explained, which is the residual error sum of squares

$$\text{RSS} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

## Results

- Hence

$$\text{TSS} = \text{RSS} + \text{ESS} .$$

- The coefficient of determination for the regression, or  $R^2$ , is

$$R^2 = \frac{\text{ESS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} .$$

- Note that TSS is purely a quantity obtained purely from the dependent variable  $Y$ .

## Degrees of Freedom

- The number of degrees of freedom of a regression with an intercept and  $p$  explanatory variables is  $k$ .
- To estimate variance of the residuals, we use the residual mean sum of squares

$$\hat{\sigma}_\epsilon^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - k}$$

## F Tests

- Suppose there are two models I and II. Model I is a subset of Model II. The excess regression sum of squares of Model II relative to Model I is

$$\begin{aligned} \text{RSS (II | I)} &= \text{RSS for Model II} - \text{RSS for Model I} \\ &= \text{ESS for Model II} - \text{ESS for Model I} \end{aligned}$$

- Correspondingly, the relative degrees of freedom of Model II with  $k_{\text{II}}$  degrees of freedom to Model I with  $k_{\text{I}}$  degrees of freedom is

$$\text{df}_{\text{II|I}} = k_{\text{II}} - k_{\text{I}}.$$

- The  $F$  statistic is

$$\frac{\text{RSS}(\text{II|I})}{\frac{k_{\text{II}} - k_{\text{I}}}{\widehat{\sigma}_{\epsilon}^2}} \sim F_{\text{df}_{k_{\text{II}} - k_{\text{I}}, n - k_{\text{II}}}}$$

## Adjusted $R^2$

- The coefficient of determination  $R^2$  is always increased by adding more explanatory variables to the linear regression model, even if they do not really explain  $Y$ .
- Recall that

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{\frac{1}{n}\text{RSS}}{\frac{1}{n}\text{TSS}}$$

where RSS is the residual sum of squares.

- The bias in  $R^2$  can be removed, resulting in the adjusted  $R^2$  denoted by  $\bar{R}^2$

$$\bar{R}^2 = 1 - \frac{\frac{1}{n-k}\text{RSS}}{\frac{1}{n-1}\text{TSS}} = 1 - \frac{\hat{\sigma}_\epsilon^2}{s^2}$$

where  $s^2$  is the sample variance of  $Y_i$ .

# Model Selection

→ Find a subset of explanatory variables that provides a parsimonious regression model by minimizing one of the following criteria

→ AIC

$$\text{AIC} = \ln(\hat{\sigma}_\epsilon^2) + \frac{2k}{n}.$$

→ BIC

$$\text{BIC} = \ln(\hat{\sigma}_\epsilon^2) + \frac{k}{n} \ln(n).$$

→  $C_p$  statistic [Mallows 1973]:

Suppose there are  $M$  explanatory variables, and  $p \leq M$ . Then  $C_p$  is defined as

$$C_p = \frac{\text{RSS}(p)}{\hat{\sigma}_{\epsilon, M}^2} - n + 2k.$$

## Newey-West Standard Errors

- \* The standard errors that are “robust” against heteroscedasticity and serial correlation were proposed by Newey and West in 1987.

- \* Define

$$\mathbf{Q} := \sum_{j=-k}^k \frac{k - |j|}{k} \frac{1}{T} \sum_{t=1}^T \hat{\epsilon}_t \mathbf{x}_t \mathbf{x}_{t-j}^\top \hat{\epsilon}_{t-j}$$

and

$$\mathbf{G} := \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$$

- \* Asymptotic variance-covariance matrix of the regression coefficients

$$\mathbb{V}(\hat{\boldsymbol{\beta}}) = \frac{1}{T} \mathbf{G}^{-1} \mathbf{Q} \mathbf{G}^{-1}$$