

Financial econometrics

Chapter 4, Multiple Linear Regression

Overview I

- 1 The multiple linear regression model
 - Multiple linear regression
 - Parametric and semi-parametric specification
- 2 Assumptions of the MLR
 - Linearity
 - Full column rank
 - Identification
 - Exogeneity
 - Spherical errors
 - Data generating process
 - Normality
- 3 The OLS estimator in MLR
 - Estimator of the regression coefficients
 - Estimator of the errors' variance
 - CAPM example
- 4 Statistical Properties of the OLS estimator
 - Unbiasedness
 - Variance of the OLS estimator
 - Finite sample distribution
 - Efficiency
 - Summary
- 5 Asymptotic Properties
 - Consistency
 - Limit distribution of the OLS
- 6 References

Section 1, The multiple linear regression model

Objectives

- Define the (multiple) linear regression model.
- Make a distinction between the semi-parametric and parametric MLR model.
- Introduce the multiple linear Gaussian model.
- Introduce a **vectorial definition** of the MLR model.

Multiple linear regression model

- Other explanatory variables might explain variations of the excess (log-) return of Intel : macroeconomic variables (e.g. inflation), financial variables (e.g. Fama-French factors).
- E.g.

$$z_{\text{intel},t} = \beta_0 + \beta_1 z_{\text{market},t} + \beta_2 \text{inflation}_t + \varepsilon_t.$$

- This is called the multiple linear regression model.

Multiple linear regression model (cont'd)

Definition (Multiple linear regression model)

The multiple linear regression model is used to study the (linear) relationship between a dependent variable and one or more independent variables. A general formulation is given by

$$y_t = \beta_0 + \beta_1 x_{t,1} + \beta_2 x_{t,2} + \cdots + \beta_K x_{t,K} + \varepsilon_t$$

where y is the dependent variable and x_1, \dots, x_K are K explanatory variables.

Notation: $x_{t,k}$ is the k -th explanatory variable for time t .

Notation

$$\underset{T \times 1}{\mathbf{y}} = (y_1, y_2, \dots, y_t, \dots, y_T)^\top$$

$$\underset{T \times 1}{\mathbf{x}_k} = (x_{1,k}, x_{2,k}, \dots, x_{t,k}, \dots, x_{T,k})^\top$$

$$\underset{T \times 1}{\boldsymbol{\varepsilon}} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_t, \dots, \varepsilon_T)^\top$$

$$\underset{T \times 1}{\boldsymbol{\beta}} = (\beta_1, \beta_2, \dots, \beta_t, \dots, \beta_T)^\top$$

Notation (cont'd)

$$\mathbf{X}_{T \times K} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots, \mathbf{x}_K)$$

or equivalently

$$\mathbf{X}_{T \times K} = \begin{pmatrix} x_{1,1} & x_{1,2} & x_{1,3} & \cdots & x_{1,k} & \cdots & x_{1,K} \\ x_{2,1} & x_{2,2} & x_{2,3} & \cdots & x_{2,k} & \cdots & x_{2,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{t,1} & x_{t,2} & x_{t,3} & \cdots & x_{t,k} & \cdots & x_{t,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{T,1} & x_{T,2} & x_{T,3} & \cdots & x_{T,k} & \cdots & x_{T,K} \end{pmatrix}$$

Multiple linear regression model (cont'd)

Definition (Multiple linear regression model)

The multiple linear regression model can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Multiple linear regression model (cont'd)

Remark: More generally, the matrix \mathbf{X} may as well contain stochastic and non stochastic elements such as:

- Constant
- Time trend
- Dummies
- etc.

Therefore, \mathbf{X} is generally a mixture of random variables and non-random variables

Multiple linear regression model (cont'd)

Remark: If the model includes a constant term (intercept), then we have

$$y_t = 1 \times \beta_1 + x_{t,2}\beta_2 + \cdots + x_{t,K}\beta_K + \varepsilon_t$$

The matrix \mathbf{X} becomes:

$$\mathbf{X}_{T \times (K+1)} = (\mathbf{1}, \mathbf{x}_2, \cdots, \mathbf{x}_K)$$

or equivalently

$$\mathbf{X}_{T \times (K+1)} = \begin{pmatrix} 1 & x_{1,2} & x_{1,3} & \cdots & x_{1,k} & \cdots & x_{1,K} \\ 1 & x_{2,2} & x_{2,3} & \cdots & x_{2,k} & \cdots & x_{2,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{t,2} & x_{t,3} & \cdots & x_{t,k} & \cdots & x_{t,K} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{T,2} & x_{T,3} & \cdots & x_{T,k} & \cdots & x_{T,K} \end{pmatrix}$$

Example

Example

The CAPM for Intel Corp. can be written as

$$z_{\text{intel},t} = \beta_0 + \beta_1 z_{\text{market},t} + \varepsilon_t,$$

or equivalently

$$\underset{T \times 1}{y} = \underset{T \times 2}{X} \underset{2 \times 1}{\beta} + \underset{T \times 1}{\varepsilon}.$$

Question: what are X , y , ε , β ?

Semi-parametric and semi-parametric specification

One key difference in the specification of the MLR:

- **Parametric model:** the distribution of the error terms is fully characterized, e.g.

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \Omega).$$

- **Semi-parametric model:** only a few moments of the error terms are specified, e.g.

$$\mathbb{E}[\varepsilon] = \mathbf{0} \quad \text{and} \quad \mathbb{V}[\varepsilon] = \mathbb{E}[\varepsilon \varepsilon^\top] = \Omega.$$

Parametric and semi-parametric specification (cont'd)

This difference does not matter for the derivation of the ordinary least square estimator. But this difference matters for (among others):

- The characterization of the statistical properties of the OLS estimator (e.g., efficiency).
- The choice of alternative estimators (e.g., the maximum likelihood estimator, etc.).

Parametric and semi-parametric specification (cont'd)

Definition (Semi-Parametric MLR)

The **semi-parametric multiple linear regression model** is defined by

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

where the error terms $\boldsymbol{\varepsilon}$ satisfies

$$\mathbb{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{0}$$

$$\mathbb{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \sigma^2 \mathbf{I}$$

and \mathbf{I} is the identity matrix (of appropriate size).

Remarks:

- If the matrix \mathbf{X} is non stochastic (fixed), i.e. there are only fixed regressors, then the conditions on the error term ε read:

$$\mathbb{E}[\varepsilon] = 0 \quad \mathbb{V}[\varepsilon] = \sigma^2 \mathbf{I}$$

- If the conditional variance-covariance matrix of ε is not diagonal, i.e. if

$$\mathbb{V}[\varepsilon|\mathbf{X}] = \Omega$$

the model is called (multiple) **generalized regression model** (GLM).

Remarks (cont'd)

Remarks:

The two conditions on the error term ε

$$\mathbb{E}[\varepsilon|\mathbf{X}] = \mathbf{0} \quad \mathbb{V}[\varepsilon|\mathbf{X}] = \sigma^2 \mathbf{I},$$

are equivalent to:

$$\mathbb{E}[\mathbf{y}|\mathbf{X}] = \mathbf{X}\beta \quad \mathbb{V}[\mathbf{y}|\mathbf{X}] = \sigma^2 \mathbf{I}.$$

Parametric and semi-parametric specification (cont'd)

Definition (multiple linear Gaussian model)

The (parametric) **multiple linear Gaussian model** is defined by

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

where the error term $\boldsymbol{\varepsilon}$ is normally distributed

$$\boldsymbol{\varepsilon}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

As a consequence, the vector \mathbf{y} has a conditional normal distribution with

$$\mathbf{y}|\mathbf{X} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

Remarks:

- The multiple linear Gaussian model is (by definition) a parametric model.
- If the matrix \mathbf{X} is non stochastic (fixed), i.e. there are only fixed regressors, then the vector \mathbf{y} has marginal normal distribution:

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

Section 2, Assumptions of the MLR

Assumptions of the MLR

The classical linear regression model consists of a set of assumptions that describes how the data set is produced by a data generating process (DGP):

- A1: Linearity
- A2: Full rank condition (or identification)
- A3: Exogeneity
- A4: Spherical errors
- A5: Data generation
- A6: Normality

Assumptions 1, linearity

Definition (Linearity)

The model is linear with respect to the parameters β_1, \dots, β_K (i.e. β).

Assumptions 1, linearity (cont'd)

Remarks:

- The models

$$y_t = \beta_0 + \beta_1 x_t + u_t$$

$$y_t = \beta_0 + \beta_1 \cos(x_t) + v_t$$

$$y_t = \beta_0 + \beta_1 \frac{1}{x_t} + \omega_t$$

are all linear w.r.t. β .

- The model

$$y_t = \beta_0 + \beta_1 x_t^{\beta_2} + \varepsilon_t$$

is not linear w.r.t. β .

- The model

$$y_t = x_t^\beta e^{\varepsilon_t}$$

can turn to linear after appropriate transformation.

Assumptions 2, full column rank

Definition (Full column rank)

\mathbf{X} is a $T \times K$ matrix with rank K .

Interpretation:

- There is no exact relationship among any of the independent variables in the model.
- The columns of \mathbf{X} are linearly **independent**.
- A matrix \mathbf{X} that is *not* full rank is also called **rank deficient**
- If the design matrix \mathbf{X} some columns (or rows) can be obtained as a linear combination of the others: if this is the case, we say that there is a (multi)-collinearity problem.
- ! Remember that a (square) matrix that is rank deficient does not have an inverse.

Assumptions 2, full column rank (cont'd)

- **Perfect** multi-collinearity (that is one variable is linearly dependent from the others) is generally not difficult to spot and is signaled by most statistical software.
- Imperfect multi-collinearity is a more serious issue.

Definition (Imperfect multicollinearity)

Imperfect multicollinearity occurs when two or more explanatory variables in a statistical model are correlated with each other, but not perfectly. I.e. they are not linearly dependent but 'almost' (e.g. one variable has ≥ 0.9 correlation with another).

Assumptions 2, Example I

Example (Multicollinearity)

Suppose that we want to estimate the following model:

$$z_{\text{intel},t} = \beta_0 + \beta_1 z_{\text{market},t} + \beta_2 (z_{\text{market},t} \times 2) + \varepsilon_t$$

The identification condition *does not hold*, $z_{\text{market},t}$ and $z_{\text{market},t} \times 2$ are perfectly collinear. It is impossible to estimate β .

Example (Full rank)

Suppose that we want to estimate the following model:

$$z_{\text{intel},t} = \beta_0 + \beta_1 z_{\text{market},t} + \beta_2 (z_{\text{market},t}^2) + \varepsilon_t$$

The identification condition *does hold*. No collinearity issues arise in estimating β .

Assumptions 2, Example II

Example (Imperfect multicollinearity)

Recall that $y = \log(1 + x) \approx x$ for $x \approx 0$. y is a non-linear transformation of x (indeed it involves the log that is non-linear), yet $y \approx x$. Look at the log-return section: by plotting x and y the relationship is pretty much linear around 0: the correlation is high.

Thus x and y are not linearly dependent (y can't be obtained from x as $a + bx$ where a, b are constants), however their correlation is very high: this is called imperfect collinearity.

Assumptions 2, Example III

Example (Questions)

$$\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \mathbf{y} = \begin{pmatrix} 0 \\ 5 \end{pmatrix}, \mathbf{z} = \begin{pmatrix} 1 \\ 2 \end{pmatrix},$$

- $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ is linearly dependent on \mathbf{x} ?
- $\begin{pmatrix} 1 \\ -8 \end{pmatrix}$ is linearly dependent on $\mathbf{x}, \mathbf{y}, \mathbf{z}$?
- $\begin{pmatrix} e^2 \\ -1 \end{pmatrix}$ is linearly dependent on $\mathbf{x}, \mathbf{y}, \mathbf{z}$?

Example (Questions)

$$\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \mathbf{y} = \begin{pmatrix} 3 \\ 3 \end{pmatrix}, \mathbf{z} = \begin{pmatrix} 10 \\ 14 \end{pmatrix},$$

Are linearly dependent?

If they are, you can write one vector as a linear combination of the others, e.g.

$$\mathbf{z} = a\mathbf{x} + b\mathbf{y}$$

Here

$$\mathbf{z} = 4a + 2b$$

But note that when taking them *two-by-two* they are linearly independent. In fact, you cannot obtain \mathbf{x} as neither $a\mathbf{y}$ not $b\mathbf{z}$. Same for \mathbf{y} and \mathbf{z} .

Linear independence I

Definition (Linear independence)

A set of vector is linearly independent if the only solution to

$$a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \cdots + a_K \mathbf{x}_K = \mathbf{0}$$

is

$$a_1 = a_2 = \cdots = a_K = 0$$

Linear independence II

Example

Form the previous example

$$2\mathbf{x} + \mathbf{y} - \frac{1}{2}\mathbf{z} = \mathbf{0}$$

We have that the null vector can be obtained as a linear combination where the coefficients are not-zero: thus $\mathbf{x}, \mathbf{y}, \mathbf{z}$ are not linearly independent.

On the other hand, we claimed that e.g., \mathbf{x} and \mathbf{y} are independent. In fact

$$a\mathbf{x} + b\mathbf{y} = \mathbf{0}$$

only if both $a = 0$ and $b = 0$.

Definition (Identification)

The multiple linear regression model is said identifiable if and only if one the following equivalent assertions holds:

- i $\text{Rank}(\mathbf{X}) = K$
- ii The matrix $\mathbf{X}^\top \mathbf{X}$ is invertible
- iii $\mathbf{X}\beta_1 = \mathbf{0}$ implies $\beta = \mathbf{0} \quad \forall \beta \in \mathbb{R}^K$
- iv $\mathbf{X}\beta_1 = \mathbf{X}\beta_2$ implies $\beta_1 = \beta_2 \quad \forall (\beta_1, \beta_2) \in \mathbb{R}^{K \times K}$

Assumptions 3, exogeneity

Strict exogeneity of the regressors

The regressors are **exogenous** if:

$$\mathbb{E}[\varepsilon|\mathbf{X}] = \mathbf{0}$$

or equivalently

$$\mathbb{E}[\varepsilon_t|x_{s,k}] = 0$$

for any explanatory variable $k \in \{1, \dots, K\}$ and any time $(t, s) \in \{1, \dots, T\}$

Assumptions 3, exogeneity (cont'd)

Remarks:

- The expected value of the error term at time t is not a function of the explanatory variables observed at any observation (including the t -th observation).
- The explanatory variables are not predictors of the error terms.
- The strict exogeneity condition can be rewritten as:

$$\mathbb{E}[\mathbf{y}|\mathbf{X}] = \mathbf{X}\beta$$

Assumptions 4, spherical errors

Spherical errors

The error terms are such that:

$$\mathbb{V}[\varepsilon_t|\mathbf{X}] = \mathbb{E}[\varepsilon_t^2|\mathbf{X}] = \sigma^2 \quad \forall t \in \{1, \dots, T\}$$

and

$$\text{Cov}(\varepsilon_t, \varepsilon_s|\mathbf{X}) = \mathbb{E}[\varepsilon_t \varepsilon_s|\mathbf{X}] = \mathbf{0} \quad \forall t \neq s$$

Notes:

- The condition of constant variances is called **homoscedasticity**.
- The uncorrelatedness across observations is called **non-autocorrelation**.

Assumptions 4, spherical errors (cont'd)

Comments:

- **Spherical disturbances** = homoscedasticity + non-autocorrelation
- If the errors are not spherical, we call them nonspherical disturbances.
- The assumption of homoscedasticity is a strong one: this is the exception rather than the rule!

Assumptions 4, spherical errors (cont'd)

Comments:

Let us consider the (conditional) variance covariance matrix of the error terms:

$$\underset{T \times T}{\mathbb{V}[\boldsymbol{\varepsilon}|\mathbf{X}]} = \mathbb{E} \left[\underset{T \times T}{\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top} | \mathbf{X} \right] =$$

$$\begin{pmatrix} \mathbb{V}[\varepsilon_1|\mathbf{X}] & \text{Cov}(\varepsilon_1\varepsilon_2|\mathbf{X}) & \cdots & \text{Cov}(\varepsilon_1\varepsilon_t|\mathbf{X}) & \cdots & \text{Cov}(\varepsilon_1\varepsilon_T|\mathbf{X}) \\ \text{Cov}(\varepsilon_2\varepsilon_1|\mathbf{X}) & \mathbb{V}[\varepsilon_2|\mathbf{X}] & \cdots & \text{Cov}(\varepsilon_2\varepsilon_t|\mathbf{X}) & \cdots & \text{Cov}(\varepsilon_2\varepsilon_T|\mathbf{X}) \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \text{Cov}(\varepsilon_t\varepsilon_1|\mathbf{X}) & \text{Cov}(\varepsilon_t\varepsilon_2|\mathbf{X}) & \cdots & \mathbb{V}[\varepsilon_t|\mathbf{X}] & \cdots & \text{Cov}(\varepsilon_t\varepsilon_T|\mathbf{X}) \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \text{Cov}(\varepsilon_T\varepsilon_1|\mathbf{X}) & \text{Cov}(\varepsilon_T\varepsilon_2|\mathbf{X}) & \cdots & \text{Cov}(\varepsilon_T\varepsilon_t|\mathbf{X}) & \cdots & \mathbb{V}[\varepsilon_T|\mathbf{X}] \end{pmatrix}$$

Assumptions 4, spherical errors (cont'd)

The assumptions of homoscedasticity and non-autocorrelation imply that:

$$\mathbb{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbb{E} \left[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top | \mathbf{X} \right] =$$
$$\begin{pmatrix} \sigma^2 & 0 & \dots & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \sigma^2 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 & \dots & \sigma^2 \end{pmatrix}$$

Assumptions 4, spherical errors (cont'd)

Notes:

- $\mathbb{V}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top|\mathbf{X}] = \sigma^2 \mathbf{I}$
 $T \times T$ $T \times T$
- homoscedasticity means the 'same variance' for all the error terms

$$\mathbb{V}[\varepsilon_1|\mathbf{X}] = \dots = \mathbb{V}[\varepsilon_T|\mathbf{X}] = \sigma^2$$

- non-autocorrelation means 'no correlation' for two error terms at two different dates

$$\text{Cov}(\varepsilon_t \varepsilon_s | \mathbf{X}) = 0 \quad \forall t \neq s$$

Assumption 5, data generating process

Data generation

The data in $(x_{t,1}, x_{t,2}, \dots, x_{t,K})$ may be any mixture of **constants** and **random variables**.

Example (Non-stochastic terms)

Some examples of non-stochastic terms used as regressors: a constant term (intercept), a time trend, or some dummy variables.

Assumption 5, data generating process (cont'd)

Comments:

- The fact that the columns of \mathbf{X} are stochastic (or not) has an impact on the asymptotic properties.
- If the explanatory variables are randomly distributed, additional assumptions regarding $(x_{t,1}, x_{t,2}, \dots, x_{t,K})$ are required. This is a statement about how the sample is drawn.
- In the sequel, we assume that $(x_{t,1}, x_{t,2}, \dots, x_{t,K})$ are **independently and identically distributed (i.i.d.)** for $t = 1, \dots, T$.

Assumption 6, normality

Normality

The data disturbances are normally distributed

$$\varepsilon|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Question: How about independence or correlations?

Assumption 6, normality (cont'd)

Comments:

- Assumption 6 implies assumption 3 (exogeneity) and 4 (spherical errors):

$$\varepsilon|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\mathbb{E}[\varepsilon|\mathbf{X}] = \mathbf{0} \quad \mathbb{V}[\varepsilon|\mathbf{X}] = \sigma^2 \mathbf{I}$$

- Normality is **not necessary** to obtain most of the results presented in the following, but practical for inference.

Section 3, The OLS estimator in MLR

Consider the MLR model

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

or equivalently, for every t

$$y_t = \sum_{k=1}^K \beta_k x_{t,k} + \varepsilon_t$$

Objective: find an estimator of $\boldsymbol{\beta}$ and σ^2 under assumptions A1-A5.

OLS estimator (cont'd)

Three equivalent approaches

- 1 Minimize the sum of squared residuals (SSR).
- 2 Use a geometric rationale/interpretation.
- 3 Solve the minimization problem with matrix notation.

Minimize the sum of squared residuals

As for the simple linear regression, we have

$$\hat{\beta} = \arg \min_{\beta} \sum_{t=1}^T \hat{\varepsilon}_t^2 = \arg \min_{\beta} \sum_{t=1}^T \left(y_t - \sum_{k=1}^K \beta_k x_{t,k} \right)^2$$

One can derive the first-order conditions with respect to β_k for $k = 1, \dots, K$ and solve a system of K equations with K unknowns.

Minimize the sum of squared residuals (cont'd)

OLS and multiple linear regression model

In the MLR model $y_t = \mathbf{x}_t^\top \boldsymbol{\beta} + \varepsilon_t$, with $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,K})^\top$, the OLS estimator $\hat{\boldsymbol{\beta}}$ is the solution of

$$\arg \min_{\boldsymbol{\beta}} \sum_{t=1}^T \left(y_t - \mathbf{x}_t^\top \boldsymbol{\beta} \right)^2$$

The **OLS estimator** of $\boldsymbol{\beta}$ is:

$$\hat{\boldsymbol{\beta}} = \left(\sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top \right)^{-1} \left(\sum_{t=1}^T \mathbf{x}_t y_t \right)$$

Geometric interpretation

- The ordinary least squares estimation methods consist in determining the adjusted vector, $\hat{\mathbf{y}}$, which is the closest to \mathbf{y} (in a certain space...) such that the squared norm between \mathbf{y} and $\hat{\mathbf{y}}$ is minimized.
- Finding $\hat{\mathbf{y}}$ is equivalent to find an estimator of β .

Geometric interpretation

The adjusted vector $\hat{\mathbf{y}}$ is the (orthogonal) projection of \mathbf{y} onto the column space of \mathbf{X} . The fitted error terms, $\hat{\epsilon}_t$, is the projection of \mathbf{y} onto the orthogonal space engendered by the column space of \mathbf{X} . The vectors $\hat{\mathbf{y}}$ and $\hat{\epsilon}_t$ are orthogonal.

Geometric interpretation (cont'd)

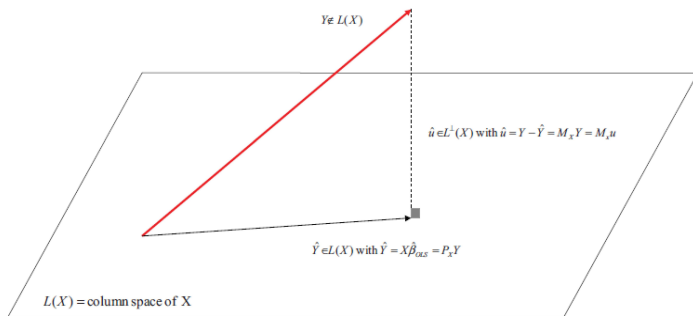


Figure: OLS as a projection.

OLS in matrix notation

OLS and multiple linear regression model

For the MLR model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ the OLS estimator $\hat{\boldsymbol{\beta}}$ is the solution of the minimization problem

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = \arg \min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

The **OLS estimator** of $\boldsymbol{\beta}$ is

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

Derivation of the OLS estimator I

In the multiple regression context, the RSS would be minimized w.r.t. all the elements in β . The RSS is the relevant loss (L) and would be given in matrix notation by

$$L = \hat{\varepsilon}^\top \hat{\varepsilon} = \hat{\varepsilon}_1^2 + \cdots + \hat{\varepsilon}_T^2 = \sum \hat{\varepsilon}_t^2.$$

Denoting the vector of estimated parameters as $\hat{\beta}$, it is possible to write

$$L = \hat{\varepsilon}^\top \hat{\varepsilon} = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\hat{\beta}^\top \mathbf{X}^\top \mathbf{y} + \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta}.$$

In order to find the parameter values that minimize loss we differentiate L w.r.t. $\hat{\beta}$ and set it to zero:

$$\frac{\partial L}{\partial \hat{\beta}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \hat{\beta} = \mathbf{0}.$$

Derivation of the OLS estimator II

Rearranging the above gives

$$\begin{aligned}2\mathbf{X}^\top \mathbf{y} &= 2\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} \\ \mathbf{X}^\top \mathbf{y} &= \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}},\end{aligned}$$

and pre-multiplying both sides by the inverse of $\mathbf{X}^\top \mathbf{X}$ gives

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

Thus the vector of OLS coefficient estimates for a set of k parameters is given by:

$$\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)^\top = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Estimate of the variance of the errors

The estimation of the variance of the errors is analogous to the SLR:

- Previously we had

$$s^2 = \frac{1}{T-2} \sum \hat{\varepsilon}_t^2$$

- Under MLR

$$s^2 = \frac{1}{T-K} \sum \hat{\varepsilon}_t^2 = \frac{\hat{\varepsilon}^\top \hat{\varepsilon}}{T-K}$$

- s^2 is unbiased for σ^2 .

CAPM example I

Example

Estimate the parameters $\beta_1, \beta_2, \sigma^2$ in the CAPM model

$$z_{\text{intel},t} = \beta_1 + \beta_2 z_{\text{market},t} \varepsilon_t$$

We consider a sample of $T = 250$ observations (1 year) from Jan. 3, 2020 to Dec. 29, 2020. We observe

$$\sum z_{\text{market},t} = 14.069 \quad \sum z_{\text{market},t}^2 = 1202.6$$

$$\sum z_{\text{intel},t} = -18.702 \quad \sum z_{\text{market},t} z_{\text{intel},t} = 1393.6$$

Question: Compute the OLS estimates of the parameters $\hat{\beta}$ and σ^2 with the matrix OLS solution.

CAPM example II

The MLR model writes as

$$\underset{250 \times 1}{\mathbf{y}} = \underset{250 \times 2}{\mathbf{X}} \underset{2 \times 1}{\boldsymbol{\beta}} + \underset{250 \times 1}{\boldsymbol{\varepsilon}}$$

Where:

- $\mathbf{y} = (z_{\text{intel},1}, z_{\text{intel},2}, \dots, z_{\text{intel},250})$
- $\underset{250 \times 2}{\mathbf{X}} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ z_{\text{market},1} & z_{\text{market},2} & \dots & z_{\text{market},250} \end{pmatrix}^{\top}$
- $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_{250})^{\top}$
- $\boldsymbol{\beta} = (\beta_1, \beta_2)^{\top}$

CAPM example III

The OLS estimator for β is:

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

with

$$\mathbf{X}_{2 \times 2}^\top \mathbf{X} = \begin{pmatrix} T & \sum z_{\text{market},t} \\ \sum z_{\text{market},t} & \sum z_{\text{market},t}^2 \end{pmatrix} = \begin{pmatrix} 250 & 14.069 \\ 14.096 & 1202.6 \end{pmatrix}$$

$$\mathbf{X}_{2 \times 1}^\top \mathbf{y} = \begin{pmatrix} \sum z_{\text{intel},t} \\ \sum z_{\text{intel},t} z_{\text{market},t} \end{pmatrix} = \begin{pmatrix} -18.702 \\ 1393.6 \end{pmatrix}$$

CAPM example IV

The OLS solution is:

$$\begin{aligned}\hat{\beta}_{2 \times 1} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \begin{pmatrix} 250 & 14.069 \\ 14.096 & 1202.6 \end{pmatrix}^{-1} \begin{pmatrix} -18.702 \\ 1393.6 \end{pmatrix} = \begin{pmatrix} -0.1401 \\ 1.1605 \end{pmatrix} \\ &= \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}\end{aligned}$$

The estimator of σ^2 is given by

$$\hat{\sigma}^2 = \frac{1}{T - K} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon} = \frac{1}{T - K} SSR$$

In this example we have $SSR = 1283.9$, then

$$\hat{\sigma}^2 = \frac{1283.9}{248} = 5.1769$$

CAPM example V

and the S.E. of the regression (RMSE) is $\sqrt{\hat{\sigma}^2} = 2.2753$.

Lastly, for variance-covariance matrix:

$$\begin{aligned}\mathbb{V}[\hat{\beta}] &= \hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})^{-1} = 5.1769 \times \begin{pmatrix} 0.00400 & -0.00004 \\ -0.00004 & 0.00083 \end{pmatrix} \\ &= \begin{pmatrix} 0.0207 & -0.0002 \\ -0.0002 & 0.0043 \end{pmatrix}\end{aligned}$$

Code example

Example

Write a code that computes the OLS estimates based on the above example.

The solution is quite simple:

```
% ... ..  
% ... ..  
% ... code to import the data and subtract Rf  
  
X      = [ones(250,1),zm];  
B      = pinv(X'*X)*X'*zi;  
Y      = X*B;  
res    = zi-Y;  
s2     = sum(res.^2)/(T-2);  
VCov   = s2*pinv(X'*X);
```

Section 4, Statistical Properties of the OLS estimator

Finite sample properties

Definition (Finite sample properties and finite sample distribution)

The finite sample properties of an estimator $\hat{\beta}$ correspond to the properties of its finite sample distribution (or exact distribution) defined for any sample size $T \in \mathbb{N}$.

Definition

Unbiased estimator Under the assumption A3 (strict exogeneity) the OLS estimator $\hat{\beta}$ is **unbiased**:

$$\mathbb{E} [\hat{\beta}] = \beta$$

where β denotes the true value of the vector of parameters. This result holds whether or not the matrix \mathbf{X} is considered as random.

Unbiasedness (cont'd)

Proof:

$$\begin{aligned}\mathbb{E}[\hat{\beta}] &= \mathbb{E}\left[\left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top \mathbf{y}\right] \\ &= \mathbb{E}\left[\left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top (\mathbf{X}\beta + \varepsilon)\right] \\ &= \left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \left(\mathbf{X}^\top \mathbf{X}\right) \beta + \left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top \mathbb{E}[\varepsilon] \\ &= \mathbf{I}\beta + \left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top \mathbb{E}[\mathbb{E}[\varepsilon|\mathbf{x}]] \\ &= \beta + 0 = \beta\end{aligned}$$

Variance of the OLS estimator

Definition

Variance of the OLS estimator Under the assumption A4 (spherical errors) the variance-covariance matrix of the OLS estimator $\hat{\beta}$ is:

$$\mathbb{V} [\hat{\beta}] = \sigma^2 \left(\mathbf{X}^\top \mathbf{X} \right)^{-1},$$

where \mathbf{X} is non-stochastic.

Variance of the OLS estimator (cont'd) I

Proof:

Remember that for a non-vector random variable x , with mean μ_x

$$\mathbb{V}[x] = \mathbb{E}[(x - \mu_x)^2]$$

If \mathbf{x} is a vector-random variable, this is similar, but in place of the square you have $\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top]$, and the variance $\mathbb{V}[\mathbf{x}]$ is actually a variance-covariance matrix: in position (i, i) of $\mathbb{V}[\mathbf{x}]$ you read the value of the variance of x_i (the i -th element of \mathbf{x}), in position (i, j) you read the covariance between x_i and x_j , $\text{Cov}(x_i, x_j)$.

Variance of the OLS estimator (cont'd) II

As above, the variance-covariance for $\hat{\beta}$ is then:

$$\mathbb{E} \left[\left(\hat{\beta} - \beta \right) \left(\hat{\beta} - \beta \right)^{\top} \right].$$

Given the OLS solution for $\hat{\beta}$ we can state that

$$\hat{\beta} = \left(\mathbf{X}^{\top} \mathbf{X} \right)^{-1} \mathbf{X}^{\top} \left(\mathbf{X} \beta + \varepsilon \right).$$

Expanding the parenthesis one gets

$$\hat{\beta} = \left(\mathbf{X}^{\top} \mathbf{X} \right)^{-1} \mathbf{X}^{\top} \mathbf{X} \beta + \left(\mathbf{X}^{\top} \mathbf{X} \right)^{-1} \mathbf{X}^{\top} \varepsilon$$

$$\hat{\beta} = \beta + \left(\mathbf{X}^{\top} \mathbf{X} \right)^{-1} \mathbf{X}^{\top} \varepsilon.$$

Variance of the OLS estimator (cont'd) III

Thus, we express the variance of $\hat{\beta}$ as

$$\begin{aligned}& \mathbb{E} \left[\left(\hat{\beta} - \beta \right) \left(\hat{\beta} - \beta \right)^{\top} \right] \\&= \mathbb{E} \left[\left(\beta + (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \varepsilon - \beta \right) \left(\beta + (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \varepsilon - \beta \right)^{\top} \right] \\&= \mathbb{E} \left[\left((\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \varepsilon \right) \left((\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \varepsilon \right)^{\top} \right] \\&= \mathbb{E} \left[(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \varepsilon \varepsilon^{\top} \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{-1} \right] \\&= (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbb{E} [\varepsilon \varepsilon^{\top}] \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{-1} \\&= (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{-1} \\&= \sigma^2 (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{-1} \\&= \sigma^2 (\mathbf{X}^{\top} \mathbf{X})^{-1}\end{aligned}$$

Variance of the OLS estimator (cont'd) IV

Therefore we have:

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

Remarks:

- $\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ is a variance-covariance matrix of the coefficients.

Question: how do we estimate it?

- The diagonal gives the estimated variances of the coefficients:

$$\text{se}(\hat{\beta}_k) = \sqrt{\left[\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})^{-1}\right]_{kk}} = \sqrt{\hat{\mathbb{V}}[\hat{\beta}_k]}$$

- The off-diagonal terms give the estimated covariance between the parameter estimates

Variance of the OLS estimator (cont'd)

Remark:

If the matrix \mathbf{X} is stochastic, the conditional variance covariance matrix of the OLS estimator $\hat{\beta}$ is

$$\mathbb{V} [\hat{\beta} | \mathbf{X}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}.$$

The unconditional variance covariance matrix is equal to

$$\mathbb{V} [\hat{\beta}] = \sigma^2 \mathbb{E}_{\mathbf{X}} \left[(\mathbf{X}^\top \mathbf{X})^{-1} \right].$$

where $\mathbb{E}_{\mathbf{X}}$ denotes the expectation with respect to the distribution of \mathbf{X} .

Finite sample distribution of the estimators

Theorem (Finite sample distribution of $\hat{\beta}$ and σ^2)

Under the assumption A6 (normality), the estimators $\hat{\beta}$ and σ^2 have finite sample distributions given by:

$$\hat{\beta} \sim \mathcal{N} \left(\beta, \sigma^2 \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \right)$$
$$\frac{\hat{\sigma}^2}{\sigma^2} (T - K) \sim \chi^2_{(T-K)}$$

Moreover, $\hat{\beta}$ and σ^2 are independent. This result holds whether or not the matrix \mathbf{X} is considered as random. In this last case, the distribution of $\hat{\beta}$ is conditional to \mathbf{X} .

Finite sample distribution of the estimators (cont'd)

Note: as a result,

$$\mathbb{E} [\hat{\sigma}^2] = \sigma^2, \quad \mathbb{V} [\hat{\sigma}^2] = \frac{2\sigma^4}{T-K}.$$

Proof, first part

- The first part is quite trivial. In $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ the stochastic component is entirely embedded in $\boldsymbol{\varepsilon} \sim \mathcal{N}$, thus the distribution of \mathbf{y} is Normal as well. Similarly, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ is normal as well. The mean and var-cov matrix of such distribution have already been computed and shown to be respectively $\boldsymbol{\beta}$ and $\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$.
- **Note:** Without A6 the mean and variance of the sampling distribution of $\hat{\boldsymbol{\beta}}$ are known, but now, under A6, we are able to characterize its distribution.

Good estimator

Question: what is a 'good' estimator of β ?

- The question is to know if there this estimator is preferred to other unbiased estimators, i.e.

$$\mathbb{V}[\beta_{OLS}] < \mathbb{V}[\beta_{other}]?$$

- This answer this question one has to use Cramer-Rao Bound (CRB) and study the efficiency of the estimator.
- The computation of the CRB is based on likelihood theory and as such requires assumptions about the distribution of ε , otherwise is just impossible to compute it.

Theorem (Efficiency of the Gaussian MLR model)

Under the assumption A6 (normality), the OLS estimator $\hat{\beta}$ is efficient. Its variance reaches the Cramer-Rao bound:

$$\mathbb{V} \left[\hat{\beta} \right] = CRB$$

Remark: The CRB expresses a lower bound on the variance of unbiased estimators of a fixed though unknown parameter, stating that the variance of *any* such estimator is at least as high the CRB. An unbiased estimator which achieves this lower bound is said to be efficient.

Remark: As a consequence, there is no other estimator with lower variance than the CRB. In this view, β_{OLS} is the best choice as it reaches the CRB and no other unbiased estimator can have smaller variance.

Problem:

- In a semi-parametric model (with no assumption on the distribution of ε), it is impossible to compute the CRB and to show the efficiency of the OLS estimator.
- The solution consists in introducing the concept of best linear unbiased estimator (BLUE): the Gauss-Markov theorem.

Theorem (Gauss-Markov theorem)

In the linear regression model under assumptions A1-A5, the least squares estimator β is the best linear unbiased estimator (BLUE) of whether \mathbf{X} is stochastic or nonstochastic.

Summary

Property	Assumption
$\hat{\beta}$ is unbiased	A3: Exogeneity
$\mathbb{V}[\hat{\beta}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$	A4: Sph. Errors
σ^2 is unbiased	A3 and A4
$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$	A6: Normality
$\frac{\hat{\sigma}^2}{\sigma^2} (T - K) \sim \chi^2_{(T-K)}$	A6: Normality

Summary (cont'd)

Property	Assumptions
$\hat{\beta}$ is BLUE	+A3, & +A4
$\hat{\beta}$ is efficient and BLUE	+ A3, + A4 & +A6 (Normality)

(In addition to, A1, A2, A5)

Section 5, Asymptotic Properties

Asymptotic Properties

Question: what is the behavior of the random variable $\hat{\beta}$ when the sample size tends to infinity?

Definition (Asymptotic theory)

Asymptotic or **large sample theory** consists in the study of the distribution of the estimator when the sample size is sufficiently large.

The asymptotic theory is fundamentally based on the notion of **convergence...**

Definition (Consistency)

An estimator $\hat{\beta}$ is consistent for β if

$$\text{plim } \hat{\beta} = \beta.$$

Limit distribution of the OLS

Theorem

*Under assumption A1-A5, the OLS estimator $\hat{\beta}$ is **asymptotically normally distributed***

$$\sqrt{N} \left(\hat{\beta} - \beta \right) \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \sigma^2 \mathbf{Q}^{-1} \right)$$

where

$$\mathbf{Q} = \text{plim} \frac{1}{N} \mathbf{X}^\top \mathbf{X} = \mathbb{E}_{\mathbf{X}} \left[\mathbf{x}_i^\top \mathbf{x}_i \right]$$

Equivalently, asymptotically:

$$\hat{\beta} \sim \mathcal{N} \left(\beta, \frac{\sigma^2}{N} \mathbf{Q}^{-1} \right)$$

Consistent estimation of the asymptotic V-Cov matrix

Remark:

The asymptotic variance is consistently estimated by the estimated variance matrix

$$\hat{\mathbb{V}}[\hat{\beta}] = s^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

where s^2 is consistent for σ^2 .

For example:

$$s^2 = \frac{1}{N - K} \hat{\epsilon}^\top \hat{\epsilon}$$

or

$$s^2 = \frac{1}{N} \hat{\epsilon}^\top \hat{\epsilon}$$

as asymptotically the correction for the number of estimated parameters is irrelevant.

Section 6, References

Disclaimer:

- Some slides from Christophe Hurlin's (University of Orleans), financial econometrics course (2019), available online.
- Some slides original (made ad-hoc for this course).
- Some slides from Colin Cameron's (University of California, Davis) lecture notes.
- Some slides based on (Hayashi, 2011, Ch. 1).

[Hay11] F. Hayashi. *Econometrics*. Princeton University Press, 2011.
ISBN: 9781400823833. URL:
<https://books.google.dk/books?id=QyIW8WUIyzcC>.