

# Linear Regression

Yutong Wang  
London School of Economics

Oct 2023

## Contents

<b>1</b>	<b>Data Generating Process - Philosophy</b>	<b>2</b>
<b>2</b>	<b>Simple Linear Regression</b>	<b>3</b>
2.1	Assumptions	3
2.1.1	Fixed Design	3
2.1.2	Random Design*	3
2.2	Sample Regression Function - Estimation	4
2.3	Memorise the OLS Solution	5
2.3.1	Estimation Strategies Based on ‘Plug-in’ Principle	5
2.3.2	OLS as ‘Plug-in’ Estimator	5
2.4	A Bunch of Concepts	6
2.4.1	Fitted Value	6
2.4.2	Sample Regression Line	6
2.4.3	Residuals	6
2.5	Sampling Distribution of OLS Estimators - Testing	7
2.5.1	The Distribution of $\hat{\beta}_1^{\text{OLS}}$	7
2.5.2	$\sigma^2$ is Known	9
2.5.3	$\sigma^2$ is Unknown - ‘Plug-in’	9
2.5.4	The Distribution of $\hat{\beta}_0^{\text{OLS}}$	10
2.5.5	$\sigma^2$ is Known	11
2.5.6	$\sigma^2$ is Unknown	11
2.6	Standard Normal, $\chi$ -squared and t-distribution*	11
2.6.1	Standard Normal	11
2.6.2	$\chi$ -squared Distribution	12
2.6.3	t-distribution	13
2.7	‘Plug-in’ Estimators*	14
2.7.1	Consistency*	14
2.7.2	Distribution*	15
<b>3</b>	<b>Multivariate Linear Regression</b>	<b>16</b>
3.1	DGP in Matrix Form	16
3.2	OLS Estimation	16
3.3	Fitted Value	17
3.4	Residual	17
3.5	Inner Product Space and Orthogonal Projection	17
3.6	$X$ is not Full-rank - towards High-dimensional Model	18
3.7	Variance (matrix calculation)	18

# 1 Data Generating Process - Philosophy

One might have seen this formula thousands of times, since it is the simplest machine-learning model, as well as a widely-used methodology in many subjects that claim they are 'quantitative'. However, it is not as easy to understand it as applying it with a few clicks in E-views/Stata/R/Python/MATLAB, etc.

Here I take the 1-dimensional linear regression model to introduce the technique to deal with it, as well as the underlying philosophy.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i. \quad (1.1)$$

PopReg

One could treat this population regression function as a **data generating process** (DGP). Data generating process could be interpreted as a process that super-nature powers (for brevity, I use 'God' to represent it.<sup>1</sup>) use 'magic' ways (or whatever similar in one's mind) to determine the state of the world.

For instance, in a specific data generating process 1.1, to determine  $y_i$ , the God has four things to do:

1. Find what is the value of **parameters** (here the parameter is  $(\beta_0, \beta_1)$ ) for this specific situation in his/her/their<sup>2</sup> dictionary;
2. Observe object  $i$ , and obtain the value of  $x_i$ ;
3. Draw a  $\varepsilon_i$  from a distribution (which is similar to generating random numbers in our computer, however, the difference is also apparent: we do not know his/her random seed, and even his/her random number generating algorithm);
4. According to the formula 1.1, add  $\beta_0$ ,  $\beta_1 x_i$  and  $\varepsilon_i$  together to determine the value of  $y_i$ .

**Remark 1.** For different problems, God may choose different parameters (different in their values as well as the dimensions) in (1) and the distribution of error terms  $\varepsilon$  in (3). For example, in the case where  $y_i$  represents the final score for students at LSE, and  $x_i$  represents how many hours they study per week in average, God chooses parameter  $(\theta_0, \theta_1)$  and distribution  $F$  to generate the data. In another case, where  $y_i$  represents whether one prefers noodles over rices,  $x_i$  represents whether one comes from the north of China, parameter  $(\alpha_0, \alpha_1)$  and distribution  $G$  are applied. We generally assume  $(\theta_0, \theta_1) \neq (\alpha_0, \alpha_1)$ , and  $F \neq G$ , which coincides with the intuition.

**Remark 2.** Here we assume that there is no measurement error, *i.e.*,  $x_i$  is the true value of the feature of the object  $i$ . This could also be understood as we put the measurement error in the  $\varepsilon$ , however, such an interpretation allows the  $\varepsilon$  to have more components, which is unnecessarily more complex.

**Remark 3.** Given the current technology level and the computing power, it is impossible to know the value of  $\beta$  and each  $\varepsilon_i$  for any  $i$ . Intuitively the  $\beta$  represents the truth of the world, which we would never know. It is somehow disappointing, however, thanks to theorems such as the Law of Large Numbers, we can approximate the truth if we have a sufficiently large amount of data. After all, in most scenarios, we do not need to be extremely accurate (except in astronautical engineering, etc.)

**Remark 4.** It is merely one way to interpret the data generating process. I introduce it since I found it compatible with and helpful in dealing with mathematical details. **I am NOT suggesting you change your beliefs.**

<sup>1</sup>This does not imply that I believe in God or not, or suggest anyone believe in it, or suggest anyone not believing in it.

<sup>2</sup>I do not specify the gender, and you could change the order for those pronouns as you like.

## 2 Simple Linear Regression

### 2.1 Assumptions

It should be noticed that, regarding whether regressors (a.k.a. features<sup>3</sup>)  $x_i$  as random, regression models can be divided two types: fixed design and random design. These two types of models naturally have different assumptions.

#### 2.1.1 Fixed Design

Under fixed design, we treat regressors  $x_i$  as fixed, which may be because we have control of  $x_i$ . For instance, in the case to find the treatment effect  $\mathbb{E}[y_i|x_i = 1]$ , where we could fully decide which individual would receive treatment and which not, implying that  $x_i$  is not random.

**Assumption 1. (Zero-mean Error)**  $\mathbb{E}[\varepsilon_i] = 0$ .

**Assumption 2. (Homoscedasticity)**  $\text{Var}(\varepsilon_i) = \sigma^2$ .

**Assumption 3. (Uncorrelated Sample)**  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \forall i \neq j$ .

**Assumption 4. (Fixed Design)**  $x_i$  is not r.v., and takes at least two values.

**Assumption 5. (Normal Error)**  $\varepsilon_i \sim N(0, \sigma^2)$ .

Note that, combining Assumption 3 and Assumption 5 gives  $\varepsilon_i \sim N(0, \sigma^2)$  i.i.d., since any uncorrelated normal r.v. are independent. In Assumption 2, ‘homoscedasticity’ means the same variance, while the counter concept is ‘heteroskedasticity’, which means the error terms have different variances.

#### 2.1.2 Random Design\*

Under random design, we treat regressors  $x_i$  as random, intuitively meaning that we can not choose on which value of  $x_i$  we observe the corresponding  $y_i$ . For example, if our goal is to analyse the effect of an individual’s income ( $x_i$ ) on the party (Demo or Repub) that he is in support for ( $y_i$ ), and we sample by randomly selecting a phone number, call and get information through the conversation. In this case, we have no control of the income ( $x_i$ ) of the person before we call him, hence should treat  $x_i$  as random.

Note that in this case,  $x_i$  is a r.v. and also has its own distribution.

**Assumption 6. (Zero-mean Error\*)**  $\mathbb{E}[\varepsilon_i|x_i] = 0$ .

**Assumption 7. (Homoscedasticity\*)**  $\text{Var}(\varepsilon_i|x_i) = \sigma^2$ .

**Assumption 8. (Uncorrelated Sample\*)**  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \forall i \neq j$ .

**Assumption 9. (Random Design\*)**  $x_i$  is not correlated to  $\varepsilon_i$  (This is actually implied by  $\mathbb{E}[\varepsilon_i|x_i] = 0$ ), and takes at least two values.

**Assumption 10. (Normal Error\*)**  $\varepsilon_i \sim N(0, \sigma^2)$ .

There is no essential difference in theory between fixed design and random design in the sense that all proofs for the random design case could be obtained from the fixed design case by replacing  $\mathbb{E}[\cdot]$  by  $\mathbb{E}[\cdot|x_i]$ . Also,

In the main part of this tutorial, we assume that the model is a fixed design for simplicity, except in chapter 2.3.

---

<sup>3</sup>‘Regressor’ is preferred in economics and econometrics, while ‘feature’ is mainly used by the machine learning community.

## 2.2 Sample Regression Function - Estimation

Once we know there is a data generating process, our interests lie in the value of parameters, and some properties of the distribution of  $\varepsilon$  (such as the variance  $\text{Var}(\varepsilon)$ <sup>4</sup>).

One of the information sources used for this purpose (estimation of parameters) that we have access to, is the data (a.k.a. ‘sample’, or ‘realisations’<sup>5</sup>). Suppose now through observation, we have data  $(x_1, y_1), \dots, (x_n, y_n)$ .

Generally, once given the data per se, and the data generating process, we have a bunch of methods to estimate the parameter. Naturally, different methods have different properties. The estimation method used in ‘Linear Regression’ is called Ordinary Least Square (OLS)<sup>6</sup>.

OLS

### Definition 1. Ordinary Least Square (OLS)

Given sample  $(x_1, y_1), \dots, (x_n, y_n)$ , and the data generating process 1.1, the OLS estimator of  $(\beta_0, \beta_1)$  is defined by:

$$\hat{\beta}^{OLS} = (\hat{\beta}_0^{OLS}, \hat{\beta}_1^{OLS}) \in \arg \min_{(\beta_0, \beta_1) \in \mathbb{R}^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2. \quad (2.2)$$

As this is an introductory level course, we will not go into details about how to evaluate an estimator, however, one should know that this is a core content in classic statistics. Jargon is provided for students who are interested in unbiased, consistent, and efficient. Also, some assumptions in the data-generating process are needed for the estimator to have these nice properties, such as exogeneity ( $\mathbb{E}[\varepsilon_i | x_i] = 0$ ).

We would give explicit expression for the solution to 1. This optimisation problem could be solved by taking derivative of  $f(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$ , since it is differentiable  $\forall (\beta_0, \beta_1) \in \mathbb{R}^2$ .

$$\begin{aligned} \frac{\partial f}{\partial \beta_0} &= \sum_{i=1}^n -2(y_i - \beta_0 - \beta_1 x_i) = 0, \\ \frac{\partial f}{\partial \beta_1} &= \sum_{i=1}^n -x_i(y_i - \beta_0 - \beta_1 x_i) = 0. \end{aligned}$$

This is a 2-dimensional linear system of  $(\beta_0, \beta_1)$ , with the solution being:

$$\beta_1^* = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2}, \quad (2.3)$$

$$\beta_0^* = \bar{y}_n - \beta_1^* \bar{x}_n. \quad (2.4)$$

**Remark 1.** Here we write the solution as  $\beta_0^*$  and  $\beta_1^*$ , since we are simply solving an optimization problem, which has nothing to do with randomness. Now let us take the randomness nature of  $(x_i, y_i)_{i=1}^n$  into consideration, set  $\hat{\beta}_0^{OLS} = \beta_0^*$ ,  $\hat{\beta}_1^{OLS} = \beta_1^*$ , we obtain the expression of OLS estimator<sup>7</sup>.

**Remark 2.** To be mathematically rigorous, one should also check the second-order derivative to show that  $(\beta_0^*, \beta_1^*)$  is the minimum point rather than maximum point.

<sup>4</sup>Note that the mean and the variance are actually defined for a distribution, rather than a random variable, since any r.v.s with the same distribution have the same mean and variance.

<sup>5</sup>My feeling: ‘sample’ is mainly used in classic statistics, while ‘data’ is used in the machine learning community, and ‘realisations’ is used in both econometrics and statistics, to emphasis both two properties of the object: as a value, and as a random variable.

<sup>6</sup>If we have assumption 5 (or 10 for random design case), the model is called parametric model, and the maximum likelihood estimation (MLE) could be applied. However, the results from MLE are the same as the one from OLS estimate.

<sup>7</sup>Recall that estimators are functions of the sample, hence random variables.

## 2.3 Memorise the OLS Solution

At the first glance, 2.3 and 2.4 seem a little bit difficult to memorise. In this section, we would provide a natural way (as a ‘plug-in’ estimator) to understand and find them.

On one hand, it could be obtained from the ‘plug-in’ principle, which is widely used in econometrics to construct new estimators. With the Continuous Mapping Theorem, we know that nice properties from the original estimators will be kept.

On the other hand, as a way to interpret the regression formula, it helps solve quick questions in linear regression (such as the question proposed by quant interviewers).

### 2.3.1 Estimation Strategies Based on ‘Plug-in’ Principle

In (1.1), on one hand, taking expectation on both sides gives:  $\mathbb{E}[y_i] = \beta_0 + \beta_1 \mathbb{E}[x_i] + \mathbb{E}[\varepsilon_i]$ . According to Assumption 1 that  $\mathbb{E}[\varepsilon_i] = 0$ , we have:

$$\mathbb{E}[y_i] = \beta_0 + \beta_1 \mathbb{E}[x_i].$$

On the other hand, notice that  $\text{Cov}(\cdot, \cdot)$  is a bi-linear function, and if we have the exogeneity assumption ( $\mathbb{E}[\varepsilon_i | x_i] = 0$ ), which implies  $\text{Cov}(x_i, \varepsilon_i) = 0$  (apply Law of Iterated Expectation to see this), then:

$$\begin{aligned} \text{Cov}(x_i, y_i) &= \text{Cov}(x_i, \beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= \text{Cov}(x_i, \beta_0) + \text{Cov}(x_i, \beta_1 x_i) + \text{Cov}(x_i, \varepsilon_i) \\ &= 0 + \beta_1 \text{Var}(x_i) + 0 \\ &= \beta_1 \text{Var}(x_i). \end{aligned}$$

Thus, in population level, we have,

$$\beta_1 = \frac{\text{Cov}(x_i, y_i)}{\text{Var}(x_i)}, \quad (2.5) \quad \text{Popu\_b1}$$

$$\beta_0 = \mathbb{E}[y_i] - \beta_1 \mathbb{E}[x_i]. \quad (2.6) \quad \text{Popu\_b0}$$

Recall the ‘plug-in’ principle: when we do not know the value of a (population level) parameter, but need it for further use, we could first estimate it, and then pretend that we ‘know’ it.

Apply ‘plug-in’ principle here, we could first estimate  $\text{Cov}(x_i, y_i)$  and  $\text{Var}(x_i)$  by  $\widehat{\text{Cov}}(x_i, y_i)$  and  $\widehat{\text{Var}}(x_i)$ , then by 2.5 we naturally construct an estimator by plugging-in  $\widehat{\text{Cov}}(x_i, y_i)$  to replace  $\text{Cov}(x_i, y_i)$ :

$$\hat{\beta}_1 = \frac{\widehat{\text{Cov}}(x_i, y_i)}{\widehat{\text{Var}}(x_i)}.$$

If further we have estimator  $\hat{\mathbb{E}}[y_i]$  and  $\hat{\mathbb{E}}[x_i]$  for parameter  $\mathbb{E}[y_i]$  and  $\mathbb{E}[x_i]$ , we could construct an estimate of  $\beta_0$  by:

$$\hat{\beta}_0 = \hat{\mathbb{E}}[y_i] - \hat{\beta}_1 \hat{\mathbb{E}}[x_i].$$

### 2.3.2 OLS as ‘Plug-in’ Estimator

Given the philosophy above, we indeed have some natural estimate of population mean, variance, and covariance:

- Sample mean (of  $x$ ):  $\hat{\mathbb{E}}[x_i] = \frac{1}{n} \sum_{i=1}^n x_i =: \bar{x}_n$ ;
- Sample mean (of  $y$ ):  $\hat{\mathbb{E}}[y_i] = \frac{1}{n} \sum_{i=1}^n y_i =: \bar{y}_n$ ;

- Sample variance:  $\widehat{\text{Var}}(x_i) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$ ;
- Sample covariance:  $\widehat{\text{Cov}}(x_i, y_i) = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_n)(x_i - \bar{x}_n)$ .

In fact, if we take these natural estimators<sup>7</sup>, and construct estimators as above, we could get the OLS estimators (Ckeck!). This tells us, the OLS estimator could be seen as a ‘plug-in’ type estimator, if we choose a proper estimate of the population mean, variance, and covariance to plug in.

I believe these natural estimators do not take time to memorise, since one should be familiar with sample mean, sample variance, sample covariance, etc. Basically, to remember the structure of OLS estimator, one needs to memorise the population version 2.6 and 2.5. The former is obtained by computing the covariance between  $x_i$  and  $y_i$ , while the latter is by taking expectation on both sides of 1.1.

## 2.4 A Bunch of Concepts

### 2.4.1 Fitted Value

By utilising the data, now we have the estimate of parameters  $\hat{\beta}_0, \hat{\beta}_1$ . With those, we could define fitted value<sup>8</sup>  $\hat{y}_i$  by:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad \forall i = 1, \dots, n. \quad (2.7)$$

SampReg

**Remark.** Fitted value could be defined for any estimated parameter (not only OLS estimate), which is why we do not use the notation  $\hat{\beta}_0^{\text{OLS}}$ .

### 2.4.2 Sample Regression Line

Geometrically, it is easy to see that these points  $(x_i, \hat{y}_i)$  are located on a straight line  $y = \hat{\beta}_0 + \hat{\beta}_1 x$ . This line is also known as the sample regression line.

### 2.4.3 Residuals

For  $i = 1, \dots, n$ , once we have the fitted value, the residual  $\hat{\varepsilon}_i$  could be defined as:

$$\begin{aligned} \hat{\varepsilon}_i &= y_i - \hat{y}_i \\ &= (\beta_0 + \beta_1 x_i + \varepsilon_i) - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \\ &= (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) x_i + \varepsilon_i. \end{aligned} \quad (2.8)$$

For interpretation in practice, *i.e.*, the viewpoint that we want to use  $x$  to explain  $y$ , residuals represent the part of the information that cannot be explained by the model ( $x$ ). From a statistics point of view,  $\hat{\varepsilon}_i$  is an estimate of  $\varepsilon_i$ , since the latter is not observable. Moreover, we calculate the arithmetic average

<sup>7</sup>Divided by  $(n-1)$  rather than  $n$ .

<sup>8</sup>In fact, for any estimate of parameter, in any data generating process, we could define fitted value. Here we only consider the OLS estimator. In general setting, we use  $f(x)$  to denote the data generating process, and  $\hat{f}(x)$  to represent the estimator.

of residuals which will be used later:

$$\begin{aligned}
\overline{\hat{\varepsilon}_n} &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \\
&= \frac{1}{n} \sum_{i=1}^n \left( (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \varepsilon_i \right) \\
&= (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)\overline{x_n} + \overline{\varepsilon_n} \\
&= \left( \beta_0 - (\overline{y_n} - \hat{\beta}_1\overline{x_n}) \right) + (\beta_1 - \hat{\beta}_1)\overline{x_n} + \overline{\varepsilon_n} \\
&= \left( \beta_0 - (\beta_0 + \beta_1\overline{x_n} + \overline{\varepsilon_n} - \hat{\beta}_1\overline{x_n}) \right) + (\beta_1 - \hat{\beta}_1)\overline{x_n} + \overline{\varepsilon_n} \\
&= 0.
\end{aligned} \tag{2.9}$$

One should be able to distinguish  $\hat{\varepsilon}_i$  and  $\varepsilon_i$ . The first one is an estimate of the second. As for their arithmetic average<sup>8</sup>, note that  $\overline{\hat{\varepsilon}_n} = 0$  for sure, but  $\overline{\varepsilon_n} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$  is a random variable, and takes 0 with 0 probability generally.<sup>9</sup>

## 2.5 Sampling Distribution of OLS Estimators - Testing

We have seen that  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are random variables (since they depend on the sample, which is a bunch of independent random variables). Thus, we could analyse the distribution of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , so that hypothesis testing tasks can be done. For example, we are always interested in  $\mathbb{H}_0 : \beta_1 = 0$ , in plain English, we hope to test whether parameters are significantly not equal to zero, which means whether there is a relationship between the two variables.

However, one should be aware that, in most cases, it is hard to derive the distribution of an estimator theoretically. Here we could do it because the error term  $\varepsilon_i$  follows a normal distribution. The famous normal distribution has been thoroughly studied and has lots of nice properties.

Without the normality of  $\varepsilon_i$ , since the estimator is a linear combination of **independent** r.v.s which is similar to ‘sample mean’, we could apply the Central Limit Theorem to draw a conclusion on the **asymptotic** distribution where we assume sample size  $n$  is large enough.

### 2.5.1 The Distribution of $\hat{\beta}_1^{\text{OLS}}$

Now let us derive the sampling distribution of OLS estimators. From 2.3, we do a bit algebra:

$$\begin{aligned}
\hat{\beta}_1 &= \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \overline{y_n})(x_i - \overline{x_n})}{\frac{1}{n} \sum_{i=1}^n (x_i - \overline{x_n})^2} \\
&= \frac{\sum_{i=1}^n (\beta_1(x_i - \overline{x_n}) + (\varepsilon_i - \overline{\varepsilon_n}))(x_i - \overline{x_n})}{\sum_{i=1}^n (x_i - \overline{x_n})^2} \\
&= \beta_1 + \frac{\sum_{i=1}^n (\varepsilon_i - \overline{\varepsilon_n})(x_i - \overline{x_n})}{\sum_{i=1}^n (x_i - \overline{x_n})^2} \\
&= \beta_1 + \frac{\sum_{i=1}^n \varepsilon_i(x_i - \overline{x_n})}{\sum_{i=1}^n (x_i - \overline{x_n})^2}.
\end{aligned} \tag{2.10}$$

In the first step, by the definition of the model 1.1, we calculate  $\overline{y_n} = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) = \beta_0 + \beta_1 \overline{x_n} + \overline{\varepsilon_n}$ , thus,  $y_i - \overline{y_n} = \beta_1(x_i - \overline{x_n}) + (\varepsilon_i - \overline{\varepsilon_n})$ . The second step is merely simple algebra.

<sup>8</sup>One should fully understand ‘mean’, ‘expectation’ and ‘(arithmetic) average’

<sup>9</sup>In fact,  $P(\overline{\varepsilon_n} \neq 0) = 1$ , since  $\overline{\varepsilon_n} \sim N(0, \frac{\sigma^2}{n})$  follows an absolutely continuous probability distribution, which takes any value  $x$  with probability 0.

However, the third step is a trick in algebra worth mentioning, since some tutorials in linear algebra use  $\frac{1}{n} \sum_{i=1}^n \varepsilon_i = 0$ , which is both wrong<sup>9</sup> and not necessary.

**Lemma 1.** *For two sequences of real number  $\{a_i\}_{i=1}^m$ ,  $\{b_i\}_{i=1}^m$  with the same length  $m$ , defining arithmetic average  $\bar{a}_m = \frac{1}{m} \sum_{i=1}^m a_i$ , and  $\bar{b}_m = \frac{1}{m} \sum_{i=1}^m b_i$  then we have:*

$$\sum_{i=1}^m (a_i - \bar{a}_m)(b_i - \bar{b}_m) = \sum_{i=1}^m a_i(b_i - \bar{b}_m) = \sum_{i=1}^m (a_i - \bar{a}_m)b_i.$$

Note that we use the notation  $a, b, m$  to emphasize that this is a basic algebra result, which does not require any assumptions related to randomness.

Now come back to 2.10. Under fixed design,  $x_i$  is not random. We could see the structure more clearly by defining the weight  $w_i = \frac{\frac{1}{n}(x_i - \bar{x}_n)}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2}$  for  $i = 1, \dots, n$ , then those weights satisfies  $\sum_{i=1}^n w_i = 0$ .

Now let us look at the second term in 2.10,

$$\hat{\beta}_1 = \beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n \varepsilon_i (x_i - \bar{x}_n)}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2} = \beta_1 + \sum_{i=1}^n w_i \varepsilon_i, \quad (2.11) \quad \boxed{\text{hatbeta1}}$$

that is, a weighted sum of i.i.d. normal (since  $\varepsilon_i \sim_{i.i.d.} N(0, \sigma^2)$ ). Thus, applying the basic property of normal distribution<sup>10</sup>,

**Proposition 1.** *If  $X \sim N(0, \sigma_1^2)$ ,  $Y \sim N(0, \sigma_2^2)$ ,  $X$  and  $Y$  are independent, then for any constant  $a, b \in \mathbb{R}$ , we have  $aX_1 + bX_2 \sim N(0, a^2\sigma_1^2 + b^2\sigma_2^2)$ .*

We could get:

$$\hat{\beta}_1 - \beta_1 = \sum_{i=1}^n w_i \varepsilon_i \sim N\left(0, \sum_{i=1}^n w_i^2 \sigma^2\right). \quad (2.12) \quad \boxed{\text{dist\_beta1}}$$

Substitute the definition of  $w_i$  in, and simple algebra gives:

$$\begin{aligned} \sum_{i=1}^n w_i^2 &= \sum_{i=1}^n \left( \frac{\frac{1}{n}(x_i - \bar{x}_n)}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2} \right)^2 \\ &= \sum_{i=1}^n \frac{\frac{1}{n^2}(x_i - \bar{x}_n)^2}{\left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right)^2} \\ &= \frac{\frac{1}{n^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2}{\left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right)^2} \\ &= \frac{\frac{1}{n}}{\left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right)} \\ &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}. \end{aligned} \quad (2.13)$$

Thus,

$$\hat{\beta}_1^{\text{OLS}} \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}\right). \quad (2.14)$$

<sup>9</sup>Keep it in mind that  $\varepsilon_i$ 's are random variables

<sup>10</sup>This property is implied by a more basic result: the summation of independent normal r.v.s' are still normal. This is one way to characterise normal distribution. Think of it: for example, if we add two uniform distribution, the result is no longer uniform.



**Remark 1.** It is also implied that  $\hat{\beta}_1^{\text{OLS}}$  is an unbiased estimator of  $\beta_1$ :

$$\mathbb{E}\hat{\beta}_1^{\text{OLS}} = \beta_1.$$

**Remark 2.** For interpretation, recall that a small variance of an estimator implied that it is a more accurate estimation <sup>11</sup>. In 2.14, we could see that, if the  $\sigma^2 = \text{Var}(\varepsilon_i)$  is smaller, or the sample size  $n$  is large, or there is more dispersion of observations around the sample mean  $\bar{x}_n$ , the estimator can be more accurate.

**Remark 3.** The algebra holds for both cases, where we know the true value  $\sigma^2$ , or we do not. However, if we want to standardise the distribution of  $\hat{\beta}_1$ , we need to consider these two cases separately.

### 2.5.2 $\sigma^2$ is Known

When  $\sigma^2 = \text{Var}(\varepsilon_i)$  is Known, 2.14 tell us:

$$\frac{\hat{\beta}_1^{\text{OLS}} - \beta_1}{\sigma} \sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sim N(0, 1). \quad (2.15)$$

If we want to conduct a hypothesis testing:  $\mathbb{H}_0 : \beta_1 = 0$  versus  $\mathbb{H}_1 : \beta_1 \neq 0$ . Under the null,  $\beta_1 = 0$ , thus,  $\frac{\hat{\beta}_1^{\text{OLS}}}{\sigma} \sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sim N(0, 1)$ . Compute the value of this test statistic, and compare it with the critical value  $z_{\alpha/2}$  (from standard normal). The task is done.

### 2.5.3 $\sigma^2$ is Unknown - ‘Plug-in’

The problem caused by unknown  $\sigma^2$  is that, we cannot compute the value of the test statistic as above. Applying ‘plug-in’ principle, we first estimate  $\sigma^2 = \text{Var}(\varepsilon_i)$  by sample variance  $\frac{1}{n-1} \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon}_n)^2$ . However,  $\varepsilon_i$  is not observable, or we say,  $\varepsilon_i$  is also unknown, so applying the ‘plug-in’ principle again, we estimate  $\varepsilon_i$  by the residual  $\hat{\varepsilon}_i$ . Eventually, we use sample variance of residuals  $\{\hat{\varepsilon}_i\}_{i=1}^n$  to estimate  $\sigma^2$ . Recall 2.11 and 2.9, do some algebra:

$$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (\hat{\varepsilon}_i - \bar{\hat{\varepsilon}}_n)^2 \quad (2.16)$$

$$\begin{aligned} &= \sum_{i=1}^n \left( \left[ (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)x_i + \varepsilon_i \right] - \left[ (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1)\bar{x}_n + \bar{\varepsilon}_n \right] \right)^2 \\ &= \sum_{i=1}^n \left( (\hat{\beta}_1 - \beta_1)(x_i - \bar{x}_n) - (\varepsilon_i - \bar{\varepsilon}_n) \right)^2 \\ &= \sum_{i=1}^n \left( (\hat{\beta}_1 - \beta_1)^2 (x_i - \bar{x}_n)^2 - 2(\hat{\beta}_1 - \beta_1)(x_i - \bar{x}_n)(\varepsilon_i - \bar{\varepsilon}_n) + (\varepsilon_i - \bar{\varepsilon}_n)^2 \right) \\ &= (\hat{\beta}_1 - \beta_1)^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 - 2(\hat{\beta}_1 - \beta_1) \sum_{i=1}^n (x_i - \bar{x}_n)(\varepsilon_i - \bar{\varepsilon}_n) + \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon}_n)^2 \\ &= \left( \sum_{i=1}^n w_i \varepsilon_i \right)^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 - 2 \left( \sum_{i=1}^n w_i \varepsilon_i \right) \sum_{i=1}^n \varepsilon_i (x_i - \bar{x}_n) + \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon}_n)^2 \\ &= - \left( \sum_{i=1}^n w_i \varepsilon_i \right)^2 \sum_{i=1}^n (x_i - \bar{x}_n)^2 + \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon}_n)^2. \end{aligned} \quad (2.17)$$

<sup>11</sup>For those interested, check out Cramer-Rao Lower Bound.

It is easy to show that  $\sigma^{-2} \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 \sim \chi^2(n-1)$ , and  $\sigma^{-2} \left( \sum_{i=1}^n w_i \varepsilon_i \right)^2 \sum_{i=1}^n (x_i - \bar{x})^2 \sim \chi^2(1)$  because of 2.12. Multiply the normal random vector  $(\varepsilon_1, \dots, \varepsilon_n)^\top$  with a proper orthonormal matrix  $O \in \mathbb{R}^{n \times n}$ , we have  $\sigma^{-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \sim \chi^2(n-2)$ , thus,  $\mathbb{E} \left[ \sigma^{-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \right] = n-2$ . By definition,  $\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$  is an unbiased estimator of  $\sigma^2$ .

Thus, replace  $\sigma^2$  by  $\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$ , we have:

$$\frac{\hat{\beta}_1^{\text{OLS}} - \beta_1}{\sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2}} \sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sim t_{n-2}. \quad (2.18)$$

If we want to conduct the hypothesis testing:  $\mathbb{H}_0 : \beta_1 = 0$  versus  $\mathbb{H}_1 : \beta_1 \neq 0$ . Under the null,  $\beta_1 = 0$ , then  $\frac{\hat{\beta}_1^{\text{OLS}}}{\sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2}} \sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \sim t_{n-2}$ . Compute the value of the test statistic, and compare it with the critical value (in  $t$ -distribution). Done.<sup>12</sup>

#### 2.5.4 The Distribution of $\hat{\beta}_0^{\text{OLS}}$

Recall that  $\hat{\beta}_0^{\text{OLS}} = \bar{y}_n - \hat{\beta}_1^{\text{OLS}} \bar{x}_n$ , combined with 2.11, we have:

$$\begin{aligned} \hat{\beta}_0^{\text{OLS}} - \beta_0 &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) - \hat{\beta}_1^{\text{OLS}} \bar{x}_n - \beta_0 \\ &= (\beta_1 - \hat{\beta}_1^{\text{OLS}}) \bar{x}_n + \bar{\varepsilon}_n \\ &= - \left( \sum_{i=1}^n w_i \varepsilon_i \right) \bar{x}_n + \bar{\varepsilon}_n \\ &= \sum_{i=1}^n \varepsilon_i \left( \frac{1}{n} - \bar{x}_n w_i \right) \\ &=: \sum_{i=1}^n \varepsilon_i v_i. \end{aligned} \quad (2.19)$$

For these weights  $v_i$ , simple algebra gives  $\sum_{i=1}^n v_i = \sum_{i=1}^n \left( \frac{1}{n} - \bar{x}_n w_i \right) = 1 - \bar{x}_n \sum_{i=1}^n w_i = 1 - \bar{x}_n \frac{\sum_{i=1}^n (x_i - \bar{x}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} = 1$ . Since  $\sum_{i=1}^n \varepsilon_i v_i$  is a linear combination of independent normal, we have:

$$\sum_{i=1}^n \varepsilon_i v_i \sim N \left( 0, \sigma^2 \sum_{i=1}^n v_i^2 \right).$$

---

<sup>12</sup>Since  $t$ -distribution has a fatter tail than standard normal, we need a larger sample size to achieve the same power. This coincides with the intuition: More information makes the estimation more accurate.

Now solve for  $\sum_{i=1}^n v_i^2$ :

$$\sum_{i=1}^n v_i^2 = \sum_{i=1}^n \left( \frac{1}{n} - \bar{x}_n w_i \right)^2 \quad (2.20)$$

$$= \sum_{i=1}^n \left( \frac{1}{n^2} - \frac{2}{n} \bar{x}_n w_i + (\bar{x}_n)^2 w_i^2 \right) \quad (2.21)$$

$$= \frac{1}{n} - 0 + (\bar{x}_n)^2 \sum_{i=1}^n w_i^2 \quad (2.22)$$

$$= \frac{1}{n} + \frac{(\bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}. \quad (2.23)$$

Thus,

$$\hat{\beta}_0^{\text{OLS}} \sim N\left(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{(\bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \right)\right) \quad (2.24)$$

### 2.5.5 $\sigma^2$ is Known

$$\frac{\hat{\beta}_0^{\text{OLS}} - \beta_0}{\sqrt{\sigma^2 \left( \frac{1}{n} + \frac{(\bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \right)}} \sim N(0, 1). \quad (2.25)$$

Under the null, Compute the value of this test statistic.

### 2.5.6 $\sigma^2$ is Unknown

With ‘plug-in’ principle, replace  $\sigma^2$  by its estimate  $\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$ , similarly:

$$\frac{\hat{\beta}_0^{\text{OLS}} - \beta_0}{\sqrt{\frac{1}{n-2} \left( \sum_{i=1}^n \hat{\varepsilon}_i^2 \right) \left( \frac{1}{n} + \frac{(\bar{x}_n)^2}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \right)}} \sim t_{n-2}. \quad (2.26)$$

## 2.6 Standard Normal, $\chi$ -squared and t-distribution\*

We have seen that in different cases, we could do testing for  $\beta$  using either  $t$ -distribution or the standard normal, and do testing for  $\sigma^2 = \text{Var}(\varepsilon)$  by using  $\chi$ -squared distribution. Here are some details about these distributions for those interested.

I leave another simple but important question here: why knowing the distribution is essential for testing? Can we do testing without knowing the distribution?

### 2.6.1 Standard Normal

It is definitely the most well-known distribution. We say  $X$  follows the standard normal distribution (denote as  $X \sim N(0, 1)$ ), if:

$$P(X \leq x) = \int_{-\infty}^x \phi(t) dt =: \Phi(x),$$

where  $\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$ . If you just finished your first-year undergrad study, and know a probability distribution should satisfy  $P(X \in \mathbb{R}) = \int_{-\infty}^{\infty} \phi(t) dt = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = 1$ , you might be interested why

this happens.

In fact, we do not have a neat expression for  $\int_{-\infty}^x e^{-t^2} dt$ . However, we have  $\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{2\pi}$ , from a quite interesting proof:

*Proof.* Let us look at the square of it:

$$\begin{aligned} \left( \int_{-\infty}^{\infty} e^{-t^2} dt \right)^2 &= \left( \int_{-\infty}^{\infty} e^{-x^2} dx \right) \left( \int_{-\infty}^{\infty} e^{-y^2} dy \right) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dx dy. \end{aligned}$$

Since the integral converges (as  $t \rightarrow \infty$ ,  $e^{-t^2}$  goes to 0 more quickly than  $\frac{1}{t^2}$ , which converges), any parametrisation will yield the same result of the integral. Set  $x = r \cos(\theta)$ ,  $y = r \sin(\theta)$ , for  $r \in [0, \infty)$  and  $\theta \in [0, 2\pi)$ , then:

$$\frac{\partial(x, y)}{\partial(r, \theta)} = \det \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{pmatrix} = \det \begin{pmatrix} \cos(\theta) & -r \sin(\theta) \\ \sin(\theta) & r \cos(\theta) \end{pmatrix} = r.$$

Then, together with  $x^2 + y^2 = r^2$ , we could continue:

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dx dy &= \int_0^{2\pi} \int_0^{\infty} e^{-r^2} r dr d\theta = 2\pi \int_0^{\infty} e^{-r^2} r dr \\ &= 2\pi \int_0^{\infty} e^{-r^2} \frac{1}{2} d(r^2) \\ &= \pi \int_0^{\infty} e^{-s} ds = \pi \end{aligned}$$

Thus,  $\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{\pi}$ , and  $\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} dt = \sqrt{2} \int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} d\frac{t}{\sqrt{2}} = \sqrt{2\pi}$ .  $\square$

The reason why normal distribution is so popular is due to the Central Limit Theorem. Basically, the theorem tells us that the standardised sample mean will be more and more like a standard normal.

**Theorem 1. Central Limit Theorem.** Suppose  $X$  follows some distribution with finite variance. Denote  $\mu = \mathbb{E}X$  and  $\sigma^2 = \text{Var}(X)$ . If  $X_1, X_2, \dots$  are independent copies of  $X$ , then we have:

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1).$$

The convergence is in the sense of distribution, i.e., For any  $t \in \mathbb{R}$ , as  $n \rightarrow \infty$ ,

$$P\left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq t\right) \longrightarrow \Phi(t).$$

It is an amazing theorem in that, we do not restrict the distribution of  $X$  too much: as long as the variance exists (i.e., the distribution is not too fat-tailed), the sample mean will eventually (as we have enough many observations) behave like normal.

### 2.6.2 $\chi$ -squared Distribution

Suppose  $X_1, X_2, \dots$  are a sequence of independent standard normal r.v.s', then we define  $\chi$ -squared distribution with degree-of-freedom  $n$  to be the distribution of

$$Y_1 := \sum_{i=1}^n X_i^2, \tag{2.27}$$

denoted as  $Y_1 \sim \chi^2(n)$  (or  $Y_1 \sim \chi_n^2$  in some books).

Due to the explicit relationship with standard normal, we get the mean and the variance:  $\mathbb{E}Y = \mathbb{E}\sum_{i=1}^n X_i^2 = \sum_{i=1}^n \mathbb{E}X_i^2 = n$ .  $\text{Var}(Y) = \text{Var}(\sum_{i=1}^n X_i^2) = \sum_{i=1}^n \text{Var}(X_i^2) = \sum_{i=1}^n \mathbb{E}(X_i^4) - [\mathbb{E}(X_i^2)]^2 = n\mathbb{E}(X_1^4) - n = 2n$ .

In the last step, we directly compute the integral or apply moment generating function method, we get  $\mathbb{E}(X_1^4) = 3n$  for  $X_1 \sim N(0, 1)$ .

To derive the probability density function (pdf) of  $\chi$ -squared distribution is a nice practice to test one's grasp of probability theory and basic calculus. Take  $n = 1$  as an example, for any  $y \geq 0$ , the cumulative probability function (cdf) is:

$$P(Y_1 \leq y) = P(-\sqrt{y} \leq X_1 \leq \sqrt{y}) = \int_{-\sqrt{y}}^{\sqrt{y}} \phi(t)dt = \int_{-\sqrt{y}}^{\sqrt{y}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2} dt.$$

Hence, the pdf is provided by taking the derivative:

$$p_{Y_1}(y) = \frac{\partial}{\partial y} \int_{-\sqrt{y}}^{\sqrt{y}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2} dt = \frac{1}{\sqrt{2\pi}} y^{-\frac{1}{2}} e^{-\frac{y}{2}}.$$

For  $n = 2$ ,

$$\begin{aligned} P(Y_2 \leq y) &= P(X_1^2 + X_2^2 \leq y) = \iint_{x_1^2 + x_2^2 \leq y} p_{X_1, X_2}(x_1, x_2) dx_1 dx_2 \\ &= \iint_{x_1^2 + x_2^2 \leq y} \phi(x_1)\phi(x_2) dx_1 dx_2 \\ &= \iint_{x_1^2 + x_2^2 \leq y} \frac{1}{2\pi} e^{-\frac{1}{2}(x_1^2 + x_2^2)} dx_1 dx_2 \\ &= \int_0^{2\pi} \int_0^{\sqrt{y}} \frac{1}{2\pi} e^{-\frac{1}{2}r^2} r dr d\theta \\ &= 2\pi \int_0^{\sqrt{y}} \frac{1}{2\pi} e^{-\frac{1}{2}r^2} d\left(\frac{r^2}{2}\right) \\ &= \int_0^{\frac{y}{2}} e^{-s} ds \\ &= 1 - e^{-\frac{y}{2}}, \end{aligned}$$

Thus,

$$p_{Y_2}(y) = \frac{\partial}{\partial y} (1 - e^{-\frac{y}{2}}) = \frac{1}{2} y e^{-\frac{y}{2}}.$$

### 2.6.3 t-distribution

Still, suppose  $X_1, X_2, \dots$  are a sequence of independent standard normal r.v.s', then we define t-distribution with degree-of-freedom  $n$  to be the distribution of

$$Z_n := \frac{X_{n+1}}{\sqrt{\sum_{i=1}^n X_i^2/n}}, \quad (2.28)$$

denoted as  $Z_n \sim t(n)$  (or  $Z_n \sim t_n$  in some books).

For t-distribution, there is an interesting result which helps foster our intuition:

**Proposition 2.** *With  $n$  goes to infinity,  $t$ -distribution with degree of freedom  $n$  converges in distribution to standard normal. In notation, suppose  $Z_n$  is defined by 2.28 for each  $n$ , then for any  $t \in \mathbb{R}$ , as  $n \rightarrow \infty$ ,*

$$P(Z_n \leq t) \rightarrow \Phi(t).$$

The Proposition can be proved by Slutsky's Theorem.

**Theorem 2. Slutsky's Theorem.**

## 2.7 'Plug-in' Estimators\*

Here is a distraction. I would like to provide a bit explanation of how and why 'plug-in' estimators work. Remember that when a parameter is unknown, and we still need to use its value, we apply the 'plug-in' principle - estimate it, and use it as the true value as we know it. This principle works fine in practice, however, in theory, some follow-up questions require further treatment.

To generally explain what are these questions, consider the following scenario: we have sample  $\{(x_i, y_i)\}_{i=1}^n$ , and want to use the sample to estimate a parameter  $\beta$  in a data generating process  $y_i = f_{\beta, \theta}(x_i) + \varepsilon_i$  (There is also another parameter  $\theta$  in this data generating process).

We construct an estimator  $\hat{\beta}$ , which requires the value of  $\theta$ . Denote the estimator as  $\hat{\beta} = g(X, Y, \theta)$  for some function  $g$ , where  $X = (x_1, \dots, x_n)^\top$ ,  $Y = (y_1, \dots, y_n)^\top$ . In the case where we know  $\theta$ , the estimator works fine.

In the case where we don't know  $\theta$ , we could construct an estimator  $\hat{\theta} = h(X, Y)$ , then plug it in  $g$  to form an estimator of  $\beta$ :  $\hat{\hat{\beta}} = g(X, Y, \hat{\theta}) = g(X, Y, h(X, Y)) =: \ell(X, Y)$  which does not require the value of  $\theta$  any more.

Obviously,  $\hat{\beta}$  and  $\hat{\hat{\beta}}$  are two different estimators, and we are interested in to what extent they are different. To assess estimators, we need some criteria. There are two widely used criteria of estimators.

### 2.7.1 Consistency\*

**Definition 2. Consistency.** *For parameter  $\beta$ , an estimator  $\hat{\beta}$  is called consistent in estimating  $\beta$  if  $\hat{\beta} \xrightarrow{p} \beta$ .*

The convergence is in the sense of probability, defined as follows.

**Definition 3. Converge in Probability (to a real value)** *We say a sequence of random variables  $\{X_n\}_{n=1}^\infty$  converge in probability to a real value  $c$  (denote as  $X_n \xrightarrow{p} c$ ), if  $\forall \delta > 0$ ,*

$$P(|X_n - c| > \delta) \longrightarrow 0, \text{ as } n \rightarrow \infty.$$

**Theorem 3. Continuous Mapping Theorem.** *If  $X_n \xrightarrow{p} c$ , and  $g(\cdot)$  is a continuous function, then  $g(X_n) \xrightarrow{p} g(c)$ .*

With the Continuous Mapping Theorem, we could prove that if both  $\hat{\theta}$  and  $\hat{\beta}$  are consistent for parameters  $\theta$  and  $\beta$ , then the 'plug-in' estimator  $\hat{\hat{\beta}}$  is consistent for  $\beta$ . This gives the desired property of a 'plug-in' estimator.

The convergence rate (how fast the convergence is) of estimators is a key sub-field in modern statistics, but we will not go into depth here. To have some kind of intuition, it would be interesting to look at the convergence rate from the information theory point of view. In the first case where we know the value of  $\theta$ , we have more information, so it is natural to expect a faster convergence.

### 2.7.2 Distribution\*

Generally, replacing a real number with a random variable will change the distribution of that estimator (r.v.). In theory, we could go through and look at how to solve for the distribution of  $g(X_1, X_2)$  where  $X_1$  and  $X_2$  are both random variables. However, the integral is typically not basic. For intuition, we give several examples:

1. Suppose  $X_1 \sim N(\mu, 1)$ . Then  $\mu X_1 \sim N(\mu^2, \mu^2)$ . Since  $\mathbb{E}X_1 = \mu$ ,  $X_1$  could be treated as an estimator of  $\mu$ . Replace the  $\mu$  in  $\mu X_1$  we get  $X_1^2$ . The point here is that  $X_1^2$  is no longer following  $N(\mu, 1)$ .
2. Given  $X \sim N(\mu, \sigma^2)$ , and i.i.d. sample  $X_1, \dots, X_n$ .  
On one hand, it is clear that  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \frac{\sigma^2}{n})$ , thus,

$$\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1).$$

On the other hand, one can use sample variance  $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  as a natural unbiased estimator of  $\sigma^2$ . If we use  $\hat{\sigma}^2$  to replace  $\sigma^2$  above, we would get:

$$\frac{\bar{X}_n - \mu}{\sqrt{\hat{\sigma}^2/n}} \sim t_{n-1}, \tag{2.29}$$

*i.e.*, the distribution changes if we replace a real value  $\sigma^2$  by a random variable  $\hat{\sigma}^2$ .<sup>10</sup> The proof of result 2.29 takes several pages, for example, see [this link](#) provided by Duke University.

---

<sup>10</sup>The result 2.29 could follow from Basu's theorem, which involves more concepts, such as 'sufficient statistic', 'ancillary statistic', and 'conditional independent'. Terms are left here for those interested.

### 3 Multivariate Linear Regression

Compared to the simple case where we only have one independent variable (a.k.a. regressors, predictors, covariates, etc.) to explain the dependent variable, here we allow the number of independent variables to be larger than 1. The population regression function becomes:

$$y_i = x_i' \theta + \varepsilon_i,$$

where  $x_i = (1, x_i^1, \dots, x_i^p)' \in \mathbb{R}^{p+1}$ ,  $\theta = (\theta_0, \theta_1, \dots, \theta_p)' \in \mathbb{R}^{p+1}$ ,<sup>13</sup> i.e., the number of independent variables is  $p$ .

#### 3.1 DGP in Matrix Form

We need to specify the sample on which we base our estimation to talk about estimation parameters  $\theta \in \mathbb{R}^{p+1}$ . As in the simple case, we have sample  $\{(x_i', y_i)\}_{i=1}^n$ . Here is the place to introduce matrix notation which helps to simplify the formula and calculation. We could incorporate all information from the sample in a single line:

$$Y = X\theta + \varepsilon,$$

where  $Y = (y_1, \dots, y_n)' \in \mathbb{R}^n$ ,  $X = (x_1, \dots, x_n)' \in \mathbb{R}^{n \times (p+1)}$ , and  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)' \in \mathbb{R}^n$ . Apply the definition of the product of matrix, it is explicitly equivalent to these  $n$  number of equations:

$$\begin{aligned} y_1 &= x_1' \theta + \varepsilon_1 \\ y_2 &= x_2' \theta + \varepsilon_2 \\ &\dots \\ y_n &= x_n' \theta + \varepsilon_n. \end{aligned}$$

#### 3.2 OLS Estimation

OLS estimator is still defined by the value of  $\theta^*$  that minimises the Residual sum of squares:

$$\theta^{\text{OLS}} \in \arg \min_{\theta} (Y - X\theta)^\top (Y - X\theta),$$

Denote  $g(\theta) = (Y - X\theta)^\top (Y - X\theta)$ , differentiate  $g$  w.r.t.  $\theta$  gives:

$$\begin{aligned} \frac{\partial g(\theta)}{\partial \theta} &= -2X^\top (Y - X\theta), \\ \frac{\partial^2 g(\theta)}{\partial \theta \partial \theta^\top} &= 2X^\top X. \end{aligned}$$

Set  $\frac{\partial g(\theta)}{\partial \theta} = 0$  yields:<sup>14</sup>

$$\hat{\theta}^{\text{OLS}} = (X'X)^{-1} X'Y,$$

and  $\frac{\partial^2 g(\theta)}{\partial \theta \partial \theta^\top} = 2X^\top X$  being positive definite ensures that  $\hat{\theta}^{\text{OLS}}$  is the minimum point. When dealing with the matrix product, always remember to check the dimension to see whether it is compatible.

<sup>13</sup>Normally, by  $\mathbb{R}^p$ , we mean  $\mathbb{R}^{p \times 1}$ , i.e., we treat vectors as column vectors by default. This is the convention in both economics and machine learning communities.

<sup>14</sup>If one is doing machine learning or statistics, it is also important to know how we get it (matrix calculus), towards Lasso regression, Ridge regression, etc.



Clearly,

$$\begin{aligned}\hat{\theta}^{\text{OLS}} &= (X'X)^{-1}X'(X\theta + \varepsilon) \\ &= \theta + (X'X)^{-1}X'\varepsilon.\end{aligned}$$

So it is the same to prove the consistency property of the estimator under assumptions  $\mathbb{E}[\varepsilon] = 0$  for fixed design (or  $\mathbb{E}[\varepsilon|X] = 0$  for random design).

### 3.3 Fitted Value

$$\begin{aligned}\hat{Y} &:= X\hat{\theta}^{\text{OLS}} \\ &= X(X'X)^{-1}X'Y \\ &= (X(X'X)^{-1}X')Y \\ &=: P_X Y.\end{aligned}$$

where  $P_X = X(X'X)^{-1}X' \in \mathbb{R}^{n \times n}$  is a symmetric idempotent matrix satisfying  $P_X' = P_X$  and  $P_X^2 = P_X$  (**check!**). From matrix theory, we know that, any eigenvalue  $\lambda$  of  $P_X$  satisfies the same equation with the matrix, i.e.  $\lambda^2 = \lambda$ , implying  $\lambda = 0, 1$ .

### 3.4 Residual

$$\begin{aligned}\hat{\varepsilon} &:= Y - \hat{Y} \\ &= Y - X\hat{\theta}^{\text{OLS}} \\ &= (I_n - X(X'X)^{-1}X')Y \\ &= (I_n - P_X)Y.\end{aligned}$$

One could see that  $(I_n - P_X)^2 = I_n - P_X + P_X - (P_X)^2 = I_n - (P_X)^2 = I_n - P_X$ , thus,  $(I_n - P_X)$  is also a projection.

### 3.5 Inner Product Space and Orthogonal Projection

For general inner product space... (some explanation later)

The vector space  $\mathbb{R}^n$  is a natural inner product space with inner product  $\langle x, y \rangle = x'y = \sum_{i=1}^n x_i y_i$ . If two vectors  $x, y$  satisfy  $\langle x, y \rangle = 0$ , we say they are orthogonal.

Coming back to the linear regression, an important result is that **the fitted value and the residual are orthogonal**:

$$\begin{aligned}\langle \hat{Y}, \hat{\varepsilon} \rangle &= (P_X Y)'(I_n - P_X)Y \\ &= Y'P_X'(I_n - P_X)Y \\ &= Y'P_X(I_n - P_X)Y \\ &= Y'OY \\ &= 0.\end{aligned}$$

This is to say,  $P_X$  is an orthogonal projection from  $\mathbb{R}^n$  onto its  $p$ -dimensional subspace spanned by  $p$  number of column vectors in  $X$ .

### 3.6 $X$ is not Full-rank - towards High-dimensional Model

It could be the case that  $X$  is not full-rank when the columns of  $X$  are linearly dependent. A simple example is that,  $x^2 = 4x^1$ , where  $x^i = (x_1^i, x_2^i, \dots, x_{n-1}^i, x_n^i)^\top$  is the  $i$ -th column of matrix  $X$ .

Then,  $X'X$  is singular ( $\text{rank}(X'X) < p + 1$ ), resulting in ill-defined  $(X'X)^{-1}$  as well as  $\hat{\theta}^{\text{OLS}}$ . The  $\hat{\theta}^{\text{OLS}}$  does not uniquely exist and can be easily seen. In fact, suppose  $\hat{\theta}^{\text{OLS}} = (\hat{\theta}_0^{\text{OLS}}, \hat{\theta}_1^{\text{OLS}}, \dots, \hat{\theta}_n^{\text{OLS}})$  minimise the  $g(\theta)$ , then we could conclude that, for any  $m \in \mathbb{R}$ , another estimator  $\check{\theta} := (\hat{\theta}_0^{\text{OLS}}, \hat{\theta}_1^{\text{OLS}} - 4m, \hat{\theta}_2^{\text{OLS}} + m, \dots, \hat{\theta}_n^{\text{OLS}})$  also minimises the  $g(\theta)$ .

However, the fitted value  $\hat{Y} = X\hat{\theta}$  is still the projection of  $Y$  on the column space of  $X$ . The difference is that the column space of  $X$  is not  $(p + 1)$ -dimensional, hence does not have a unique way to be represented. In practice, it is implied that there are redundant features (regressors), and we should remove them first before doing the regression.

It is common to see non-full-rank  $X$  in lots of applications, such as image analysis. In this case,  $p$  is even larger than  $n$ . The typical solution is to add more assumptions to the data-generating process, such as the sparsity: among all  $p$  regressors, only  $s$  (where  $s \ll p$ ) number of them are functioning (only  $s$  number of  $\theta_i, i = 0, \dots, p$  are not 0), but we don't know which. To simultaneously estimate which regressor is functioning and its coefficient, researchers have found models with regularization (especially Lasso regression) work perfectly. We will talk about it later (if I have time lol).

### 3.7 Variance (matrix calculation)

You will find it easier to calculate in matrix form than in scalar form since lots of stuff just cancel out.

$$\begin{aligned} \text{Var}(\hat{\theta}|X) &= \text{Var}\left((X'X)^{-1}X'Y|X\right) \\ &= \left((X'X)^{-1}X'\right) \cdot \text{Var}(Y|X) \cdot \left((X'X)^{-1}X'\right)' \\ &= \left((X'X)^{-1}X'\right) \cdot (\sigma^2 I_n) \cdot \left((X'X)^{-1}X'\right)' \\ &= \sigma^2 \left((X'X)^{-1}X'\right) \left((X'X)^{-1}X'\right)' \\ &= \sigma^2 \left((X'X)^{-1}X'\right) X(X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1}, \end{aligned}$$

where  $\text{Var}(Y|X) = \sigma^2 I_n$  is because  $\text{Var}(y_i|x_i) = \sigma^2$ . Here we have used the associativity of matrix multiplication and commutativity between  $I$  and any matrices.

### 3.8 Multivariate Normal (Gaussian)

Similar to uni-variate normal, we could define multi-variate normal as<sup>15</sup>:

$$X \sim N(\mu, \Sigma), \quad (3.30)$$

where  $X = (X_1, \dots, X_p)' \in \mathbb{R}^p$  is the normal random vector,  $\mu = (\mu_1, \dots, \mu_p)' \in \mathbb{R}^p$  and  $\mu_i = \mathbb{E}X_i$ .  $\Sigma \in \mathbb{R}^{p \times p}$  is the covariance matrix,  $\Sigma_{ij} = \text{Cov}(X_i, X_j)$ , thus,  $\Sigma$  is symmetric since  $\text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i)$ .

For the random vector defined in 3.30, the joint probability density function is: <sup>16</sup>

$$\phi(x) = (2\pi)^{-p/2} (\det \Sigma)^{-1/2} \exp \left\{ -\frac{1}{2} (x - \mu)' \Sigma^{-1} (x - \mu) \right\} \quad (3.31)$$

<sup>15</sup>In some books, it is denoted as  $X \sim MVN(\mu, \Sigma)$ .

<sup>16</sup>The fact that  $\Sigma$  is non-singular coincides with our intuition: If  $\Sigma$  is singular, then there exist one of the  $X_i$  being the linear combination of the other, which means we should not embed the vector  $X$  into  $p$ -dimension.

For multivariate normal, a nice property is important:

**Proposition 3.** *If  $X \sim N(\mu, \Sigma)$  as in 3.30, and  $S \in \mathbb{R}^{m \times p}$  is any  $m \times p$  matrix ( $m \in \mathbb{Z}_+$ ). Then,*

$$SX \sim N(S\mu, S\Sigma S').$$

*Proof.* The proposition follows from three results:

1. Any linear combination of normal (not necessary to be independent) is still normal.
2. For any  $a \in \mathbb{R}^p$ ,  $\mathbb{E}a'X = a'\mathbb{E}X$ .
3.  $\text{Cov}(a'Y, b'Z) = \text{Cov}\left(\sum_{i=1}^p a_i Y_i, \sum_{i=1}^p b_i Z_i\right) = \sum_{i=1}^p \sum_{i=1}^p a_i b_i \text{Cov}(Y_i, Z_i)$ .

By applying the first result, we know each coordinate of  $SX$  follows a univariate normal distribution, thus, we only need to make sure what's their mean/variance/covariance. For mean, use the second, for variance and covariance, use the third.  $\square$

With this result, we could draw conclusions about the distribution of  $\hat{\theta}$ .

$$\hat{\theta} \sim N\left(\theta, \sigma^2(X'X)^{-1}\right)$$