

---

SESSION 1 : Exercises

---

## 1 Unidimensional model

**Assumptions and notations** The data is generated according to this model.

$$\begin{aligned} y_i &= \theta_0^* + \theta_1^* x_i + \varepsilon_i, \\ \varepsilon_i &\stackrel{i.i.d}{\sim} \varepsilon, \text{ for } i = 1, \dots, n \\ \mathbb{E}(\varepsilon) &= 0 \end{aligned}$$

The statistical problem is to estimate the parameters  $\hat{\theta}_0$  and  $\hat{\theta}_1$  from the observations of  $y$  and  $x$ . with the *ordinary least squares* (OLS) method.

**EXERCICE 1.** Give the estimator for the coefficient for the basic model.

**Solution.** The OLS method consists in minimizing the sum of squared residuals, i.e., the differences between observed values and predicted values :

$$\min \sum_{i=1}^n \varepsilon_i^2 = \min \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \min \sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i)^2 \equiv \min f(\hat{\theta}_0, \hat{\theta}_1),$$

Check the first order condition

$$\frac{\partial f}{\partial \hat{\theta}_0} = 0, \quad \frac{\partial f}{\partial \hat{\theta}_1} = 0.$$

From the first equation :

$$\frac{\partial f}{\partial \hat{\theta}_0} = 0 \quad \Rightarrow \quad -2 \sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) = 0 \quad \Rightarrow \quad \boxed{\hat{\theta}_0 = \bar{y} - \bar{x} \hat{\theta}_1}.$$

From the second equation,

$$\begin{aligned} \frac{\partial f}{\partial \hat{\theta}_1} = 0 &\iff -2 \sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) x_i = 0 \\ &\iff \sum_{i=1}^n (y_i x_i - \bar{y} x_i - \hat{\theta}_1 x_i^2 + \hat{\theta}_1 x_i^2) = 0 \\ &\iff \boxed{\hat{\theta}_1 = \frac{\sum_{i=1}^n y_i x_i - n \bar{y} \bar{x}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2}} \end{aligned}$$

The Hessian of  $f$  is positive definite iff  $\text{Var}(x) > 0$ , so this critical point is indeed a minimum iff  $\text{Var}(x) > 0$ .

$$f(\hat{\theta}_0, \hat{\theta}_1) = \sum_{i=1}^n (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i)^2,$$

$$\text{Hessian } H = \begin{bmatrix} \frac{\partial^2 f}{\partial \hat{\theta}_0^2} & \frac{\partial^2 f}{\partial \hat{\theta}_0 \partial \hat{\theta}_1} \\ \frac{\partial^2 f}{\partial \hat{\theta}_1 \partial \hat{\theta}_0} & \frac{\partial^2 f}{\partial \hat{\theta}_1^2} \end{bmatrix} = \begin{bmatrix} 2n & 2 \sum_{i=1}^n x_i \\ 2 \sum_{i=1}^n x_i & 2 \sum_{i=1}^n x_i^2 \end{bmatrix},$$

$$\det(H) = 4 \left( n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2 \right).$$

which is greater than 0 iff  $\text{Var}(x) > 0$ .

**EXERCICE 2.** In the slides we give the estimator for the coefficient for the basic model when the data is centered.

**Solution.** Defining centered variables :

$$x_i^c = x_i - \bar{x}, \quad y_i^c = y_i - \bar{y}, \quad \hat{y}_i^c = \hat{y}_i - \bar{y}, \quad \hat{\varepsilon}_i = y_i^c - \hat{y}_i^c,$$

the OLS method can be rewritten as :

$$y_i^c = \hat{\theta}_0 + \hat{\theta}_1 x_i^c + \hat{\varepsilon}_i.$$

We want to minimize the sum of squared residuals using centered data :  $\sum_{i=1}^n \hat{\varepsilon}_i^2$ .

$$\frac{\partial f}{\partial \hat{\theta}_0^c} = -2 \sum_{i=1}^n (y_i^c - \hat{\theta}_0^c - \hat{\theta}_1 x_i^c) = 0 \quad \Rightarrow \quad \hat{\theta}_0^c = \bar{y}^c - \hat{\theta}_1 \bar{x}^c = 0,$$

$$\frac{\partial f}{\partial \hat{\theta}_1} = -2 \sum_{i=1}^n (y_i^c - \hat{\theta}_0^c - \hat{\theta}_1 x_i^c) x_i^c = 0 \quad \Rightarrow \quad \hat{\theta}_1^c = \frac{\sum_{i=1}^n x_i^c y_i^c}{\sum_{i=1}^n (x_i^c)^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \hat{\theta}_1.$$

Since  $\hat{\theta}_0^c = 0$ , the regression model with centered data simplifies to :

$$\hat{y}_i^c = \hat{\theta}_1 x_i^c$$

**EXERCICE 3.** In the slides we derive a state the relationship between the slope  $\hat{\theta}_1$  and the correlation  $\text{corr}_n(\mathbf{y}, \mathbf{x})$

**Solution.** The linear correlation coefficient  $\text{corr}_n(\mathbf{y}, \mathbf{x})$  measures the degree of linear covariation between  $y$  and  $x$  and is defined as :

$$\text{corr}_n(\mathbf{y}, \mathbf{x}) = \frac{\text{cov}(x, y)}{\sqrt{\text{Var}(x)\text{Var}(y)}} = \frac{\frac{1}{n} \sum_{i=1}^n y_i x_i - \bar{y}\bar{x}}{\sqrt{\left(\frac{1}{n} \sum_{i=1}^n y_i^2 - \bar{y}^2\right) \left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2\right)}} = \hat{\theta}_1 \frac{\sqrt{\text{Var}(x)}}{\sqrt{\text{Var}(y)}}$$

where  $\hat{\theta}_1 =$ .

We can check these properties

- 1)  $\text{corr}_n(\mathbf{y}, \mathbf{x}) \in [-1; 1]$  :
- 2)  $\text{corr}_n(\mathbf{y}, \mathbf{x})$  is dimensionless.
- 3)  $\text{corr}_n(\mathbf{y}, \mathbf{x})$  is symmetric :  $\text{corr}_n(\mathbf{y}, \mathbf{x}) = \text{corr}_n(\mathbf{x}, \mathbf{y})$ .
- 4)  $\text{corr}_n(\mathbf{y}, \mathbf{x})$  is not affected by a change of variable :

**EXERCICE 4.** Show that the variance decomposes as

$$\sum_{i=1}^n (y_i^c)^2 = \sum_{i=1}^n (\hat{y}_i^c)^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2.$$

**Solution.** Analysis of variance allows decomposing the total variance into explained variance and residual variance in order to measure the quality of the regression model. By definition, we have  $\hat{\varepsilon}_i = y_i - \hat{y}_i \iff y_i^c = \hat{\varepsilon}_i + \hat{y}_i^c$  with  $\hat{y}_i^c = \hat{\theta}_1 x_i^c$ , hence :

$$\sum_{i=1}^n (y_i^c)^2 = \sum_{i=1}^n (\hat{y}_i^c)^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 + 2 \sum_{i=1}^n \hat{y}_i^c \hat{\varepsilon}_i.$$

Let's see what happens with  $\sum_{i=1}^n \hat{y}_i^c \hat{\varepsilon}_i$

$$\begin{aligned} \sum_{i=1}^n \hat{\varepsilon}_i \hat{y}_i^c &= \sum_{i=1}^n (y_i^c - \hat{y}_i^c) \hat{y}_i^c = \sum_{i=1}^n (y_i^c - \hat{y}_i^c) \hat{\theta}_1 x_i^c \\ &= \hat{\theta}_1 \left( \sum_{i=1}^n (y_i^c - \hat{\theta}_1 x_i^c) x_i^c \right) = \hat{\theta}_1 \left( \sum_{i=1}^n y_i^c x_i^c - \hat{\theta}_1 \sum_{i=1}^n (x_i^c)^2 \right). \end{aligned}$$

Also, plug-in the slope in the previous equation

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n x_i^c y_i^c}{\sum_{i=1}^n (x_i^c)^2} \implies \sum_{i=1}^n x_i^c y_i^c - \hat{\theta}_1 n \mathbb{V}\text{ar}(x) = 0 \implies \sum_{i=1}^n \hat{\varepsilon}_i \hat{y}_i^c = 0.$$

Therefore

$$\sum_{i=1}^n (y_i^c)^2 = \sum_{i=1}^n (\hat{y}_i^c)^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2.$$

This is the analysis of variance equation that describes the decomposition of the total variability of the point cloud into explained variations and residual variations. Indeed :

$$\begin{aligned} \text{SST} &= \sum_{i=1}^n (y_i^c)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (\text{variance of } y, \text{ sum of squares total}) \\ \text{SSR} &= \sum_{i=1}^n (\hat{y}_i^c)^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (\text{variance of } \hat{y}, \text{ sum of squares regression}) \\ \text{SSE} &= \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\text{variance of } \varepsilon, \text{ sum of squares error}), \end{aligned}$$

with SST the total sum of squares, SSE the explained sum of squares (by the regression line), and SSR the residual sum of squares. We write the analysis of variance :

$$\text{SST} = \text{SSE} + \text{SSR}.$$

**EXERCICE 5.** Show that the  $R^2$ , defined as the ratio of the SSR and the SST, is equivalent to

$$1 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (y_i^c)^2}.$$

**Solution.** The ordinary least squares estimates for simple linear regression are

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}, \quad \hat{\theta}_1 = \frac{\text{cov}(x, y)}{\mathbb{V}\text{ar}(x)}.$$

The coefficient of determination  $R^2$  for simple linear regression is defined as

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n (\hat{\theta}_0 + \hat{\theta}_1 x_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Substituting  $\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$  yields

$$R^2 = \frac{\sum_{i=1}^n (\hat{\theta}_1 (x_i - \bar{x}))^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

This simplifies to

$$R^2 = \hat{\theta}_1^2 \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2} = \hat{\theta}_1^2 \frac{\text{Var}(x)}{\text{Var}(y)}.$$

Since

$$\hat{\theta}_1 = \frac{\text{cov}(x, y)}{\text{Var}(x)},$$

then

$$R^2 = \left( \frac{\text{cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}} \right)^2.$$

So we conclude

$$\boxed{R^2 = \text{corr}(x, y)^2}.$$

## 2 Multidimensional case

**Assumptions and notations** We assume that the data is generated following the model :  $\mathbf{y} = X\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$  where :

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,p} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \theta_0^* \\ \vdots \\ \theta_p^* \end{pmatrix}}_{\boldsymbol{\theta}^*} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}}$$

and so  $y_i = \theta_0^* + \sum_{j=1}^p \theta_j^* x_{i,j} + \varepsilon_i$ . Using the matrix notation, we can notate the following quantities as :

Total variance of the data (SST) :

$$\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|_2^2 = \sum_{i=1}^n (y_i - \bar{y})^2$$

The goal of OLS is

$$\hat{\boldsymbol{\theta}} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{p+1}} \|\mathbf{y} - X\boldsymbol{\theta}\|_2^2 = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^{p+1}} \sum_{i=1}^n \left[ y_i - \left( \theta_0 + \sum_{j=1}^p \theta_j x_{i,j} \right) \right]^2.$$

The minimizer is unique iff

$$\text{Ker}(X) = \{0\}.$$

Assuming this condition, the solution is denoted by  $\hat{\boldsymbol{\theta}}$ , and the fitted values (predictions) are

$$\hat{\mathbf{y}} = X\hat{\boldsymbol{\theta}}.$$

and the residual, the error in the prediction is (SSE) :

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - X\hat{\boldsymbol{\theta}}$$

We can then decompose the variability as follows :

- Variance explained by the regression (SSR) :

$$\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|_2^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

- Variance of the error (SSE) :

$$\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Finally, the decomposition holds (proof done bellow) :

$$\underbrace{\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|_2^2}_{\text{SST}} = \underbrace{\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|_2^2}_{\text{SSR}} + \underbrace{\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}_{\text{SSE}}.$$

**EXERCICE 6.** The solution of the normal equations is unique when the columns of  $X$  are linearly independent, i.e., iff  $X$  is full rank,  $\text{rank}(X) = p+1$  iff its kernel is trivial :  $\text{Ker}(X) = \{0\}$ . We want to show here that

$$\text{Ker}(X) = \text{Ker}(X^\top X).$$

**Solution.**

- 1) First, show  $\text{Ker}(X) \subseteq \text{Ker}(X^\top X)$  : Let  $\boldsymbol{\theta} \in \text{Ker}(X)$ , then  $X\boldsymbol{\theta} = 0$ . Multiplying both sides by  $X^\top$ , we get

$$X^\top X\boldsymbol{\theta} = X^\top 0 = 0,$$

so  $\boldsymbol{\theta} \in \text{Ker}(X^\top X)$ .

- 2) Next, show  $\text{Ker}(X^\top X) \subseteq \text{Ker}(X)$  : Let  $\boldsymbol{\theta} \in \text{Ker}(X^\top X)$ , then  $X^\top X\boldsymbol{\theta} = 0$ . Consider the quadratic form

$$\boldsymbol{\theta}^\top X^\top X\boldsymbol{\theta} = \|X\boldsymbol{\theta}\|^2 = 0.$$

Hence,  $X\boldsymbol{\theta} = 0$ , so  $\boldsymbol{\theta} \in \text{Ker}(X)$ .

**EXERCICE 7.** Show that the residuals are centered,  $\sum \hat{\epsilon}_i = 0$ .

**Solution.**

Recall that the OLS estimator  $\hat{\boldsymbol{\theta}}$  satisfies the normal equations

$$X^\top X\hat{\boldsymbol{\theta}} = X^\top \mathbf{y}.$$

Equivalently,

$$X^\top (X\hat{\boldsymbol{\theta}} - \mathbf{y}) = \mathbf{0}_{p+1},$$

that is,

$$X^\top \hat{\boldsymbol{\epsilon}} = \mathbf{0}_{p+1}.$$

In words : the residual vector  $\hat{\boldsymbol{\epsilon}}$  is orthogonal to all the columns of  $X$ . Since the first column of  $X$  is the vector  $\mathbf{1}_n$  (corresponding to the intercept), we get

$$\hat{\boldsymbol{\epsilon}}^\top \mathbf{1}_n = \sum_{i=1}^n \hat{\epsilon}_i = 0.$$

**EXERCICE 8.** Show that the variance decomposes as

$$\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|^2. \quad (1)$$

**Solution.**

From the normal equations we know that the residuals are orthogonal to the fitted values :

$$\langle \hat{\epsilon}, \hat{\mathbf{y}} \rangle = 0,$$

and, since  $\mathbf{1}_n$  (the intercept column) belongs to the column space of  $X$ , we also have

$$\langle \hat{\epsilon}, \bar{\mathbf{y}}\mathbf{1}_n \rangle = 0.$$

Thus,

$$\langle \hat{\epsilon}, \hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n \rangle = 0.$$

Now write

$$\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n = (\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n).$$

Taking squared norms and using the orthogonality relation above gives

$$\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|^2.$$

This is exactly the variance decomposition formula (1).

**EXERCICE 9.** Show that  $0 \leq R^2 \leq 1$  and

$$R^2 = 1 - \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2} \quad (2)$$

**Solution.**

Recall the definition of  $R^2$

$$R^2 = \frac{\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|^2}{\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2}.$$

Reordering the Equality (1) :

$$\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|^2 = \|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2 - \|\mathbf{y} - \hat{\mathbf{y}}\|^2.$$

Dividing both sides by  $\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2$ , we obtain

$$R^2 = \frac{\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\mathbf{1}_n\|^2}{\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2} = 1 - \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\|\mathbf{y} - \bar{\mathbf{y}}\mathbf{1}_n\|^2}.$$