

Lecture 9 | 14.04.2025

# Statistical inference

in a linear model (asymptotics)

# Overview

## □ Normal linear regression model

- **Assumptions:** random sample  $(Y_i, \mathbf{X}_i^\top)^\top$  for  $i = 1, \dots, n$  from the joint distribution  $F_{(Y, \mathbf{X})}$  such that  $Y_i | \mathbf{X}_i \sim N(\mathbf{X}_i^\top \beta, \sigma^2)$
- **Inference:** confidence intervals for  $\beta_j$ , confidence regions for  $\beta$  and linear combinations of the form  $\mathbb{L}\beta$  (plus the corresponding statistical tests)

## □ Linear regression model without normality

### Assumptions (A1):

- random sample  $(Y_i, \mathbf{X}_i^\top)^\top$ ,  $i = 1, \dots, n$  from the joint distribution  $F_{(Y, \mathbf{X})}$
- mean specification  $E[Y_i | \mathbf{X}_i] = \mathbf{X}_i^\top \beta$ , respectively  $E[\mathbf{Y} | \mathbb{X}] = \mathbb{X}\beta$
- thus, for errors  $\varepsilon_i = Y_i - \mathbf{X}_i^\top \beta$  we have  $E[\varepsilon_i | \mathbf{X}_i] = E[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = 0$  and  $\text{Var}(\varepsilon_i | \mathbf{X}_i) = \text{Var}[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = \text{Var}[Y_i | \mathbf{X}_i] = \sigma^2(\mathbf{X}_i)$
- and for unconditional expectations,  $E[\varepsilon_i] = E[E[\varepsilon_i | \mathbf{X}_i]] = 0$  and  $\text{Var}(\varepsilon_i) = \text{Var}(E[\varepsilon_i | \mathbf{X}_i]) + E[\text{Var}(\varepsilon_i | \mathbf{X}_i)] = \text{Var}(0) + E[\sigma^2(\mathbf{X}_i)] = E[\sigma^2(\mathbf{X}_i)]$

### Inference:

- involves different confidence intervals statistical tests of hypotheses

# Parameter estimation without normality

- In the **normal regression model**  $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  one can simply use the distributional specification to formulate the likelihood (loglikelihood)
- In a **general regression model**  $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  where  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \boldsymbol{\Sigma})$  the likelihood (loglikelihood resp.) can not be formulated (the distribution is missing)
- The most common approach in this case is based on the method of least squares (LSE), thus, the vector of the estimated parameters is given as

$$\hat{\boldsymbol{\beta}}_n = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{Arg max}} \sum_{i=1}^n \left[ Y_i - \mathbf{x}_i^\top \boldsymbol{\beta} \right]^2$$

- and the estimated vector of parameters can be given explicitly as

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

which is the **BLUE** estimate for  $\boldsymbol{\beta} \in \mathbb{R}^p$  but for the statistical inference we need to know its (asymptotic) distributional properties (how does this random quantity behave when  $n \in \mathbb{N}$  tends to infinity,  $n \rightarrow \infty$ )

## Some additional assumptions

The random sample  $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$  drawn from some joint distribution  $F_{(Y, \mathbf{X})}$  of a generic  $(p + 1)$ -dimensional random vector  $(Y, \mathbf{X}^\top)^\top$ . Let  $\mathbf{X} = (X_1, \dots, X_p)^\top$ . Let the following holds:

### Assumptions (A2):

- ▣  $E|X_j X_k| < \infty$  for  $j, k \in \{1, \dots, p\}$
- ▣  $E(\mathbf{X}\mathbf{X}^\top) = \mathbb{W} \in \mathbb{R}^{p \times p}$  is a positive definite (regular) matrix
- ▣  $\mathbb{V} = \mathbb{W}^{-1}$

Note, that the assumptions stated above refer to the population model—the population properties

# Empirical counterparts for $\mathbb{W}$ and $\mathbb{V}$

- Both matrices,  $\mathbb{W} \in \mathbb{R}^{p \times p}$  and  $\mathbb{V} \in \mathbb{R}^{p \times p}$  are theoretical (population) characteristics, the dimensions are fixed for any  $n \in \mathbb{N}$ , and they are typically not known in practical applications
- Both matrices can be however estimated using the empirical data—the observed random sample  $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$
- Define the following:
  - $\mathbb{W}_n = \mathbb{X}^\top \mathbb{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$
  - $\mathbb{V}_n = \mathbb{W}_n^{-1}$  if it exists (eventually it will for  $n \in \mathbb{N}$  large enough)
- Under the assumptions in (A1) and (A2)
  - $\frac{1}{n} \mathbb{W}_n \longrightarrow \mathbb{W}$  a.s. (in P) as  $n \rightarrow \infty$
  - $n \mathbb{V}_n \longrightarrow \mathbb{V}$  a.s. (in P) as  $n \rightarrow \infty$

It is also good to realize that  $(\mathbb{X}^\top \mathbb{X})^{-1}$  may not exist for any  $n \in \mathbb{N}$  but as far as  $\frac{1}{n}(\mathbb{X}^\top \mathbb{X})$  converges almost surely (in probability) to the matrix  $\mathbb{W}$  (positive definite) we also have that  $P(\text{rank}(\mathbb{X}^\top \mathbb{X}) = p) \rightarrow 1$ , for  $n \rightarrow \infty$

# Problems of the statistical inference

Analogously as in the normal linear model, the statistical inference concerns confidence sets and statistical tests about  $\beta \in \mathbb{R}^p$  and its linear combinations

- Statistical inference can be performed with respect to the parameters  $\beta$  and  $\sigma^2$  but, it can be also of some interest to do inference about some (appropriate) linear combination(s) of  $\beta$
- From the practical point of view, we are interested in the parameter vector  $\beta$  itself but also linear combinations of the form  $\mathbf{I}^\top \beta$  or  $\mathbb{L}\beta$

The estimates for the unknown parameters  $\beta \in \mathbb{R}^p$  and  $\sigma^2 > 0$  are

$$\square \hat{\beta}_n = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i Y_i \right) \quad (\text{LSE})$$

$$\square s_n^2 = \frac{1}{n-p} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \frac{1}{n-p} \|\mathbf{Y} - \mathbb{X} \hat{\beta}\|_2^2, \text{ where } \hat{Y}_i = \mathbf{x}_i^\top \hat{\beta} \quad (\text{MSe})$$

Both estimates, quantities  $\hat{\beta}_n$  and  $\hat{s}_n^2$ , are random quantities (random vector and random variable) and, therefore, it is reasonable to investigate their statistical properties (e.g., mean, variance, distribution, etc.)

# Homoscedastic vs. heteroscedastic model

Recall, that in the assumption in (A1) the conditional variance of  $\varepsilon_i$  depends on  $\mathbf{X}_i$ , which is reflected by the notation  $\text{Var}(\varepsilon_i|\mathbf{X}_i) = \sigma^2(\mathbf{X}_i)$

**Two situations are typically distinguished:**

□ **Homoscedastic model)** (Assumption A3a)

$$\sigma^2(\mathbf{X}) = \text{Var}(Y|\mathbf{X}) = \sigma^2 > 0$$

□ **Heteroscedastic model** (Assumption A3b)

$\sigma^2(\mathbf{X}) = \text{Var}(Y|\mathbf{X})$  such that  $E[\sigma^2(\mathbf{X})] < \infty$  and moreover, it also holds that  $E[\sigma^2(\mathbf{X})X_jX_k] < \infty$  for  $j, k \in \{1, \dots, p\}$

# Consistency of the LSE estimates

□ In particular, we are interested in the following parameters:

□  $\beta \in \mathbb{R}^p$

□  $\sigma^2 > 0$

□  $\theta = \mathbf{I}^\top \beta \in \mathbb{R}$ , for some nonzero vector  $\mathbf{I} \in \mathbb{R}^p$

□  $\Theta = \mathbf{L}\beta \in \mathbb{R}^m$ , for some matrix  $\mathbf{L} \in \mathbb{R}^{m \times p}$  with linearly independent rows

□ The corresponding estimates are defined straightforwardly and it holds (under (A1), (A2), and (A3a/A3b)) that

□  $\hat{\beta}_n \longrightarrow \beta$  a.s. (in P), for  $n \rightarrow \infty$

□  $\hat{\theta}_n = \mathbf{I}^\top \hat{\beta}_n \longrightarrow \theta$  a.s. (in P), for  $n \rightarrow \infty$

□  $\hat{\Theta}_n = \mathbf{L}\hat{\beta}_n \longrightarrow \Theta$ , a.s. (in P), for  $n \rightarrow \infty$

□ Under the homoscedastic model ((A1), (A2), and (A3a)) it also holds

□  $\hat{s}_n^2 \longrightarrow \sigma^2$ , a.s. (in P), for  $n \rightarrow \infty$



# Asymptotic normality

Under the assumptions stated in (A1), (A2), and (A3a) and, additionally, for  $E[\varepsilon^2 X_j X_k] < \infty$  for  $j, k = 1, \dots, p$  the following holds:

- $\sqrt{n}(\hat{\beta}_n - \beta) \xrightarrow{\mathcal{D}} N_p(\beta, \sigma^2 \mathbb{V})$  for  $n \rightarrow \infty$
- $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{D}} N(0, \sigma^2 \mathbf{I}^\top \mathbb{V} \mathbf{I})$ , as  $n \rightarrow \infty$
- $\sqrt{n}(\hat{\Theta}_n - \Theta) \xrightarrow{\mathcal{D}} N_m(\mathbf{0}, \sigma^2 \mathbf{L} \mathbb{V} \mathbf{L}^\top)$ , as  $n \rightarrow \infty$

# Statistical inference based on asymptotics

- Define the random variable

$$T_n = \frac{\mathbf{I}^\top \hat{\boldsymbol{\beta}}_n - \mathbf{I}^\top \boldsymbol{\beta}}{\sqrt{MSe \cdot \mathbf{I}^\top (\mathbb{X}^\top \mathbb{X})^{-1} \mathbf{I}}} \left( = \frac{\sqrt{n}(\mathbf{I}^\top \hat{\boldsymbol{\beta}}_n - \mathbf{I}^\top \boldsymbol{\beta})}{\sqrt{\sigma^2 \mathbf{I}^\top \mathbb{V} \mathbf{I}}} \cdot \sqrt{\frac{\sigma^2 \mathbf{I}^\top \mathbb{V} \mathbf{I}}{MSe \cdot \mathbf{I}^\top \left[ n(\mathbb{X}^\top \mathbb{X})^{-1} \right] \mathbf{I}}} \right)$$

$\hookrightarrow$  where it is easy to see that the first term in the brackets converges (in distribution) to  $N(0, 1)$  and the second term converges (in probability) to one (Cramér-Slutsky)

- Define the random variable

$$Q_n = \frac{(\mathbb{L} \hat{\boldsymbol{\beta}}_n - \mathbb{L} \boldsymbol{\beta})^\top [\mathbb{L}(\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{L}^\top]^{-1} (\mathbb{L} \hat{\boldsymbol{\beta}}_n - \mathbb{L} \boldsymbol{\beta})}{MSe}$$

$\hookrightarrow$  where  $\sqrt{n}(\mathbb{L} \hat{\boldsymbol{\beta}}_n - \mathbb{L} \boldsymbol{\beta}) \xrightarrow{\mathcal{D}} N(0, \sigma^2 \mathbb{L} \mathbb{V} \mathbb{L}^\top)$  for  $n \rightarrow \infty$  and, also,  $(MSe \cdot [\mathbb{L} n(\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{L}^\top])^{-1} \xrightarrow{\mathcal{D}} \sigma^2 \mathbb{L} \mathbb{V} \mathbb{L}^\top$  for  $n \rightarrow \infty$  (Cramér-Slutsky)

- Then it holds that

- $T_n \xrightarrow{\mathcal{D}} N(0, 1)$  for  $n \rightarrow \infty$

- $Q_n \xrightarrow{\mathcal{D}} \chi_m^2$  for  $n \rightarrow \infty$

$\hookrightarrow$  what is this good for in practical applications and inference?

# Standard inference tools – summary

In general, the **statistical inference** is a (mathematical) process of using observed data (e.g., random sample) to make **valid and consistent conclusions** or predictions about an unknown (much larger) population. It involves (mainly) the hypotheses testing and confidence intervals construction.

## ❑ Confidence intervals

- ❑ normal linear regression model (exact coverage)
- ❑ linear regression model without normality (asymptotic coverage)

## ❑ Statistical tests

- ❑ normal linear regression model (based on the exact distribution)
- ❑ linear regression model without normality (asymptotic validity)