

Lecture 7 | 14.04.2026

Statistical inference

in a linear model without normality (asymptotics)

Overview

□ Normal linear regression model

- **Assumptions:** Random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ from the joint distribution $F_{(Y, \mathbf{X})}$, such that $Y_i | \mathbf{X}_i \sim N(\mathbf{X}_i^\top \boldsymbol{\beta}, \sigma^2)$ while the marginal distribution of r.v. $\mathbf{X} \in \mathbb{R}^p$ does not depend on $\boldsymbol{\beta} \in \mathbb{R}^p$ neither on $\sigma^2 > 0$
- **Inference:** Confidence intervals/tests for $\beta_j \in \mathbb{R}$, confidence regions/tests for $\boldsymbol{\beta} \in \mathbb{R}^p$ (or a subset of parameters), including also linear combinations (reparametrizations) of $\boldsymbol{\beta} \in \mathbb{R}^p$ in a form $\boldsymbol{\ell}^\top \boldsymbol{\beta} \in \mathbb{R}$ or $\mathbf{L}\boldsymbol{\beta} \in \mathbb{R}^m$

Overview

□ Normal linear regression model

- **Assumptions:** Random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ from the joint distribution $F_{(Y, \mathbf{X})}$, such that $Y_i | \mathbf{X}_i \sim N(\mathbf{X}_i^\top \beta, \sigma^2)$ while the marginal distribution of r.v. $\mathbf{X} \in \mathbb{R}^p$ does not depend on $\beta \in \mathbb{R}^p$ neither on $\sigma^2 > 0$
- **Inference:** Confidence intervals/tests for $\beta_j \in \mathbb{R}$, confidence regions/tests for $\beta \in \mathbb{R}^p$ (or a subset of parameters), including also linear combinations (reparametrizations) of $\beta \in \mathbb{R}^p$ in a form $\ell^\top \beta \in \mathbb{R}$ or $\mathbf{L}\beta \in \mathbb{R}^m$

□ 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 $\varepsilon_i = Y_i - \mathbf{X}_i^\top \beta$ it holds $E[\varepsilon_i | \mathbf{X}_i] = E[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = 0$ and $\text{Var}(\varepsilon_i | \mathbf{X}_i) = \text{Var}[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = \text{Var}[Y_i | \mathbf{X}_i] = \sigma^2(\mathbf{X}_i)$

Overview

□ Normal linear regression model

- **Assumptions:** Random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ from the joint distribution $F_{(Y, \mathbf{X})}$, such that $Y_i | \mathbf{X}_i \sim N(\mathbf{X}_i^\top \beta, \sigma^2)$ while the marginal distribution of r.v. $\mathbf{X} \in \mathbb{R}^p$ does not depend on $\beta \in \mathbb{R}^p$ neither on $\sigma^2 > 0$
- **Inference:** Confidence intervals/tests for $\beta_j \in \mathbb{R}$, confidence regions/tests for $\beta \in \mathbb{R}^p$ (or a subset of parameters), including also linear combinations (reparametrizations) of $\beta \in \mathbb{R}^p$ in a form $\ell^\top \beta \in \mathbb{R}$ or $\mathbf{L}\beta \in \mathbb{R}^m$

□ 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 $\varepsilon_i = Y_i - \mathbf{X}_i^\top \beta$ it holds $E[\varepsilon_i | \mathbf{X}_i] = E[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = 0$ and $\text{Var}(\varepsilon_i | \mathbf{X}_i) = \text{Var}[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = \text{Var}[Y_i | \mathbf{X}_i] = \sigma^2(\mathbf{X}_i)$
- Unconditional characteristics: $E[\varepsilon_i] = E[E[\varepsilon_i | \mathbf{X}_i]] = 0$ and $\text{Var}(\varepsilon_i) = \text{Var}(E[\varepsilon_i | \mathbf{X}_i]) + E[\text{Var}(\varepsilon_i | \mathbf{X}_i)] = \text{Var}(0) + E[\sigma^2(\mathbf{X}_i)] = E[\sigma^2(\mathbf{X}_i)]$

Overview

□ Normal linear regression model

- **Assumptions:** Random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ from the joint distribution $F_{(Y, \mathbf{X})}$, such that $Y_i | \mathbf{X}_i \sim N(\mathbf{X}_i^\top \beta, \sigma^2)$ while the marginal distribution of r.v. $\mathbf{X} \in \mathbb{R}^p$ does not depend on $\beta \in \mathbb{R}^p$ neither on $\sigma^2 > 0$
- **Inference:** Confidence intervals/tests for $\beta_j \in \mathbb{R}$, confidence regions/tests for $\beta \in \mathbb{R}^p$ (or a subset of parameters), including also linear combinations (reparametrizations) of $\beta \in \mathbb{R}^p$ in a form $\ell^\top \beta \in \mathbb{R}$ or $\mathbf{L}\beta \in \mathbb{R}^m$

□ 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 $\varepsilon_i = Y_i - \mathbf{X}_i^\top \beta$ it holds $E[\varepsilon_i | \mathbf{X}_i] = E[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = 0$ and $\text{Var}(\varepsilon_i | \mathbf{X}_i) = \text{Var}[Y_i - \mathbf{X}_i^\top \beta | \mathbf{X}_i] = \text{Var}[Y_i | \mathbf{X}_i] = \sigma^2(\mathbf{X}_i)$
- Unconditional characteristics: $E[\varepsilon_i] = E[E[\varepsilon_i | \mathbf{X}_i]] = 0$ 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:

- Again, involves confidence intervals/regions and statistical tests of hypotheses about the unknown parameters ($\beta \in \mathbb{R}^p$ and $\sigma^2 > 0$)

Parameter estimation without normality

- In the normal regression model $\mathbf{Y} = \mathbb{X}\beta + \varepsilon$, one can simply use the distributional specification to formulate the whole likelihood (loglikelihood)
- In a general regression model $\mathbf{Y} = \mathbb{X}\beta + \varepsilon$ where $\varepsilon \sim (\mathbf{0}, \sigma^2\mathbb{I})$, the likelihood (the loglikelihood resp.) can not be formulated as the underlying distribution is not specified

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 whole likelihood (loglikelihood)
- In a general regression model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\mathbb{I})$, the likelihood (the loglikelihood resp.) can not be formulated as the underlying distribution is not specified
- 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

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

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 whole likelihood (loglikelihood)
- ❑ In a general regression model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\mathbb{I})$, the likelihood (the loglikelihood resp.) can not be formulated as the underlying distribution is not specified
- ❑ 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

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

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 whole likelihood (loglikelihood)
- In a general regression model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2\mathbb{I})$, the likelihood (the loglikelihood resp.) can not be formulated as the underlying distribution is not specified
- 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 as the sample size increases, i.e., $n \rightarrow \infty$)

Additional assumptions for asymptotics

Let $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ be a random sample drawn from some joint distribution $F_{(Y, \mathbf{X})}$ of some generic $(p + 1)$ -dimensional random vector $(Y, \mathbf{X}^\top)^\top \in \mathbb{R}^{p+1}$. Let $\mathbf{X} = (X_1, \dots, X_p)^\top$ and 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
- $\mathbb{V} = \mathbb{W}^{-1}$

Additional assumptions for asymptotics

Let $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$ be a random sample drawn from some joint distribution $F_{(Y, \mathbf{X})}$ of some generic $(p + 1)$ -dimensional random vector $(Y, \mathbf{X}^\top)^\top \in \mathbb{R}^{p+1}$. Let $\mathbf{X} = (X_1, \dots, X_p)^\top$ and 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 all refer to the population model—the population properties related to the random vector $\mathbf{X} \in \mathbb{R}^p$.

The first two items in Assumptions (A2) are, indeed, technical conditions that need to be satisfied (for asymptotics). The last item is more like an introduction of a notation.

Empirical counterparts for matrices \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 both these matrices are typically not known in practical applications
- Both matrices can be, however, effectively estimated using the empirical data—the observed random sample $\{(Y_i, \mathbf{X}_i^T)^\top; i = 1, \dots, n\}$

Empirical counterparts for matrices \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 both these matrices are typically not known in practical applications
- Both matrices can be, however, effectively estimated using the empirical data—the observed random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$
- Define the following:
 - $\mathbb{W}_n = \mathbf{X}^\top \mathbf{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ and, also
 - $\mathbb{V}_n = \mathbb{W}_n^{-1}$ if it exists (eventually it will for $n \in \mathbb{N}$ large enough)

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

Empirical counterparts for matrices \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 both these matrices are typically not known in practical applications
- Both matrices can be, however, effectively estimated using the empirical data—the observed random sample $\{(Y_i, \mathbf{X}_i^\top)^\top; i = 1, \dots, n\}$
- Define the following:
 - $\mathbb{W}_n = \mathbb{X}^\top \mathbb{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ and, also
 - $\mathbb{V}_n = \mathbb{W}_n^{-1}$ if it exists (eventually it will for $n \in \mathbb{N}$ large enough)
- Under the assumptions in (A1) and (A2), the following hold:
 - $\frac{1}{n} \mathbb{W}_n \rightarrow \mathbb{W}$ a.s. (in P) as $n \rightarrow \infty$
 - $n \mathbb{V}_n \rightarrow \mathbb{V}$ 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 some (rather small) values of $n \in \mathbb{N}$ but as far as $\frac{1}{n}(\mathbb{X}^\top \mathbb{X})$ converges almost surely (in probability) to the matrix \mathbb{W} (which is 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 unknown parameters $\beta \in \mathbb{R}^p$ and $\sigma^2 > 0$ but, it can be also of some interest to do inference about some (appropriate) linear combination(s) of β
- From the practical point of view, it may be often of some more relevant interest to say something about a specific linear combination—for instance, $I^\top \beta \in \mathbb{R}$ or $L\beta \in \mathbb{R}^m$ rather than $\beta \in \mathbb{R}^p$ itself

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 unknown parameters $\beta \in \mathbb{R}^p$ and $\sigma^2 > 0$ but, it can be also of some interest to do inference about some (appropriate) linear combination(s) of β
- From the practical point of view, it may be often of some more relevant interest to say something about a specific linear combination—for instance, $I^\top \beta \in \mathbb{R}$ or $L\beta \in \mathbb{R}^m$ rather than $\beta \in \mathbb{R}^p$ itself

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

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

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

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 unknown parameters $\beta \in \mathbb{R}^p$ and $\sigma^2 > 0$ but, it can be also of some interest to do inference about some (appropriate) linear combination(s) of β
- From the practical point of view, it may be often of some more relevant interest to say something about a specific linear combination—for instance, $\mathbf{1}^\top \beta \in \mathbb{R}$ or $\mathbf{L}\beta \in \mathbb{R}^m$ rather than $\beta \in \mathbb{R}^p$ itself

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

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

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

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 ε_i depends on \mathbf{X}_i , which is reflected by the notation $\text{Var}(\varepsilon_i|\mathbf{X}_i) = \sigma^2(\mathbf{X}_i)$

Homoscedastic vs. heteroscedastic model

Recall, that in the assumption in (A1), the conditional variance of ε_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)

$$\sigma^2(\mathbf{X}) = \text{Var}(Y|\mathbf{X}) \text{ such that } E[\sigma^2(\mathbf{X})] < \infty$$

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

Homoscedastic vs. heteroscedastic model

Recall, that in the assumption in (A1), the conditional variance of ε_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)

$$\sigma^2(\mathbf{X}) = \text{Var}(Y|\mathbf{X}) \text{ such that } E[\sigma^2(\mathbf{X})] < \infty$$

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

The **homoscedastic model** (without the normality) can be generally expressed as

$$Y_i|\mathbf{X}_i \sim (\mathbf{X}_i^\top \beta, \sigma^2)$$

while the **heteroscedastic model** (again, without the normality assumption) as

$$Y_i|\mathbf{X}_i \sim (\mathbf{X}_i^\top \beta, \sigma^2(\mathbf{X}_i))$$

for some positive variance function $\sigma^2 : \mathbb{R}^p \rightarrow (0, \infty)$

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

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

- Under Assumptions (A1) and Assumptions (A2), homoscedastic model implies the heteroscedastic one, i.e. (A3a) \implies (A3b)

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

- Under Assumptions (A1) and Assumptions (A2), homoscedastic model implies the heteroscedastic one, i.e. (A3a) \implies (A3b)

- The corresponding estimates are defined straightforwardly and it holds (under (A1), (A2), and (A3a/A3b)) that
 - $\widehat{\beta}_n \rightarrow \beta$ a.s. (in P), for $n \rightarrow \infty$
 - $\widehat{\theta}_n = \mathbf{I}^\top \widehat{\beta}_n \rightarrow \theta$ a.s. (in P), for $n \rightarrow \infty$
 - $\widehat{\Theta}_n = \mathbf{L}\widehat{\beta}_n \rightarrow \Theta$, a.s. (in P), for $n \rightarrow \infty$

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

- Under Assumptions (A1) and Assumptions (A2), homoscedastic model implies the heteroscedastic one, i.e. (A3a) \implies (A3b)

- The corresponding estimates are defined straightforwardly and it holds (under (A1), (A2), and (A3a/A3b)) that
 - $\widehat{\beta}_n \longrightarrow \beta$ a.s. (in P), for $n \rightarrow \infty$
 - $\widehat{\theta}_n = \mathbf{I}^\top \widehat{\beta}_n \longrightarrow \theta$ a.s. (in P), for $n \rightarrow \infty$
 - $\widehat{\Theta}_n = \mathbf{L}\widehat{\beta}_n \longrightarrow \Theta$, a.s. (in P), for $n \rightarrow \infty$

- Under the homoscedastic model ((A1), (A2), and (A3a)) it also holds
 - $\widehat{s}_n^2 \longrightarrow \sigma^2$, a.s. (in P), for $n \rightarrow \infty$

Asymptotic normality (homoscedastic model)

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:

Asymptotic normality (homoscedastic model)

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:

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

Asymptotic normality (homoscedastic model)

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:

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

Analogous (asymptotic normality) result will also hold under the assumptions (A1), (A2), and (A3b) – i.e., under the heteroscedastic model – however, the asymptotic variance term will be slightly different. In both cases, the proof follows the same lines and, thus, it will be given for the homoscedastic model only.

The asymptotic normality above already provides all necessary inference tools!

Statistical inference based on asymptotics

- Define the random variable T_n as

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

↔ 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)

Statistical inference based on asymptotics

- Define the random variable T_n as

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

↔ 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 (a quadratic form) D_n as

$$Q_n = \frac{(\mathbf{L} \widehat{\boldsymbol{\beta}}_n - \mathbf{L} \boldsymbol{\beta})^\top \left[\mathbf{L} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{L}^\top \right]^{-1} (\mathbf{L} \widehat{\boldsymbol{\beta}}_n - \mathbf{L} \boldsymbol{\beta})}{MSe}$$

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

Statistical inference based on asymptotics

- Define the random variable T_n as

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

\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 (a quadratic form) D_n as

$$Q_n = \frac{(\mathbf{L} \widehat{\boldsymbol{\beta}}_n - \mathbf{L} \boldsymbol{\beta})^\top \left[\mathbf{L} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{L}^\top \right]^{-1} (\mathbf{L} \widehat{\boldsymbol{\beta}}_n - \mathbf{L} \boldsymbol{\beta})}{MSe}$$

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

- Then it holds (asymptotically) that

- $T_n \xrightarrow{\mathcal{D}} N(0, 1)$ for $n \rightarrow \infty$
- $Q_n \xrightarrow{\mathcal{D}} \chi_m^2$ for $n \rightarrow \infty$

Standard inference tools – overview

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/prediction intervals/sets construction.

Standard inference tools – overview

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/prediction intervals/sets 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)