Lecture 10 | 28.04.2025

Linear regression models with heteroscedasetic errors

Normal linear model

Assumptions

□ random sample (Y_i, X_i) for i = 1, ..., n from some joint distribution function $F_{(Y,X)}$, such that $Y_i | X_i \sim N(X_i^\top \beta, \sigma^2)$

u regression model of the form $Y_i = X_i^\top \beta + \varepsilon_i$

Inference

- □ confidence intervals for $\beta_j \in \mathbb{R}$, confidence regions for $\beta \in \mathbb{R}^p$, and linear combinations of the form $\mathbb{L}\beta$ for some $\mathbb{L} \in \mathbb{R}^{m \times p}$
- \square parameter estimates $\widehat{\beta}$ (constructed in terms of LSE or MLE) are BLUE and the follow the normal distribution

 $\widehat{\boldsymbol{\beta}} \sim N_{p}(\boldsymbol{\beta}, \sigma^{2}(\mathbb{X}^{\top}\mathbb{X})^{-1})$

The statistical inference is exact and it is based on the normal distribution (if the variance parameter is known) or the Student's *t*-distribution or Fisher's *F*-distribution respectively for $\sigma^2 > 0$ unknown

Linear model without normality

Assumptions (A1)

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

Assumptions (A2)

- $\Box E|X_jX_k| < \infty \text{ for } j,k \in \{1,\ldots,p\}$
- $\square E(\boldsymbol{X}\boldsymbol{X}^{\top}) = \mathbb{W} \in \mathbb{R}^{p \times p} \text{ is a positive definite matrix}$

$$\square \mathbb{V} = \mathbb{W}^{-1}$$

Assumptions (A3a/A3b)

- □ Homoscedastic model) $\sigma^2(\mathbf{X}) = Var(Y|\mathbf{X}) = \sigma^2 > 0$
- □ Heteroscedastic model $\sigma^2(\mathbf{X}) = Var(Y|\mathbf{X})$ such that $E[\sigma^2(\mathbf{X})] < \infty$ and moreover, it also holds that $E[\sigma^2(\mathbf{X})X_jX_k] < \infty$ for $j, k \in \{1, ..., p\}$

Inference under (A1), (A2), and (A3b)

Inference (without normality + homoscedastic errors)

- □ confidence intervals for $\beta_j \in \mathbb{R}$, confidence regions for $\beta \in \mathbb{R}^{\rho}$, and linear combinations of the form $\mathbb{L}\beta$ for some $\mathbb{L} \in \mathbb{R}^{m \times \rho}$
- □ parameter estimates $\hat{\beta}_n$ (sometimes also $\hat{\beta}$), constructed in terms of LSE or MLE, are BLUE, they are consistent (convergence in probability) and they follow asymptotically the normal distribution

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \quad \stackrel{\mathcal{D}}{\longrightarrow} \quad N_p(\mathbf{0}, \sigma^2 \mathbb{V})$$

The statistical inference is approximate/assymptocal and it is based on the normal distribution (regardless of whether the variance $\sigma^2 > 0$ is known or unknown)

Note that

$$\sqrt{n} \cdot \widehat{\beta}_n = \sqrt{n} (\mathbb{X}^\top \mathbb{X})^{-1} (\mathbb{X}^\top \underbrace{(\mathbb{X}\beta + \varepsilon)}_{\mathbf{Y}}) = \sqrt{n} \cdot \underbrace{\mathbb{V}_n \mathbb{V}_n^{-1} \beta}_{\beta} + \underbrace{n \mathbb{V}_n}_{\rightarrow \mathbb{V}} \cdot \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \varepsilon_i}_{(\star)}$$

 \hookrightarrow where (*) converges (in distribution) to $N_p(\mathbf{0}, E[\sigma^2(\mathbf{X})\mathbf{X}\mathbf{X}^\top])$ (Central Limit Theorem)

General linear model (heteroscedasticity)

u random sample (Y_i, X_i) for i = 1, ..., n from the joint distribution $F_{(Y, X)}$

- \Box mean specification $E[\mathbf{Y}|\mathbb{X}] = \mathbb{X}\beta$, for $\beta \in \mathbb{R}^{p}$
- □ variance specification $Var[\mathbf{Y}|\mathbb{X}] = \sigma^2 \mathbb{W}^{-1}$, for some known matrix $\mathbb{W} \in \mathbb{R}^{n \times n}$ (positive definite)
- generally, the normal distribution is not assumed, therefore

 $\boldsymbol{Y}|\mathbb{X} \sim (\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbb{W}^{-1})$

Example

Consider a linear regression model, where the dependent variables Y_i for i = 1, ..., n represent some averages across $w_i \in \mathbb{N}$ independent subjects, where for each subject we assume the same variance (i.e., a homoscedastic model for the subjects)

General least squares

Consider a general linear model $\mathbf{Y} | \mathbb{X} \sim (\mathbb{X}\beta, \sigma^2 \mathbb{W}^{-1})$ where $rank(\mathbb{X}) = p < n$ (where $\mathbb{X} \in \mathbb{R}^{n \times p}$). Than the following holds:

□
$$\hat{\beta} = (\mathbb{X}^{\top} \mathbb{W} \mathbb{X})^{-1} \mathbb{X}^{\top} \mathbb{W} \mathbf{Y}$$
 is BLUE for $\beta \in \mathbb{R}^{p}$
□ $\hat{\mu} = \hat{\mathbf{Y}} = \mathbb{X} \hat{\beta}$ is BLUE for $\mu = E[\mathbf{Y}|\mathbb{X}]$
□ for $\mathbf{I} \in \mathbb{R}^{p}$, where $\mathbf{I} \neq \mathbf{0}$, $\mathbf{I}^{\top} \hat{\beta}$ is BLUE for $\theta = \mathbf{I}^{\top} \beta$
□ $MSe_{\mathcal{G}} = \frac{1}{n-p} \|\mathbb{W}^{1/2}(\mathbf{Y} - \hat{\mathbf{Y}})\|_{2}^{2}$ is unbiased estimate of $\sigma^{2} > 0$

If, additionaly, $\mathbf{Y} | \mathbb{X} \sim N(\mathbb{X}\beta, \sigma^2 \mathbb{W}^{-1})$ then the estimates $\widehat{\beta} \in \mathbb{R}^p$ follow the corresponding normal distribution and, moreover,

$$\frac{MSe_G(n-p)}{\sigma^2} = \frac{SSe_G}{\sigma^2} \sim \chi^2_{n-p}$$

and SSe and $\widehat{\mathbf{Y}}$ are conditionally, given $\mathbb X,$ mutually independent

General linear model – utilization

- the general linear model is typically used with partially aggregated data—mostly in a way, that instead of raw observations we observe independent averages over specific classes (that we can control for with the set of the regressor variables)
- □ if the estimation of the mean structure is of the interest only, the aggregated data can be also replicated and the correponding mean estimates will be the same
- however, if there is also some interest in the variance estimation (e.g., there is a need to perform some statistical inference), the model based on the replicated data will fail (the variance estimates are artificially underestimated—e.g., too short confidence intervals)
- □ the situations described above all refer to a diagonal (weighting) matrix \mathbb{W} . However, in general, the matrix $\mathbb{W} \in \mathbb{R}^{n \times n}$ can have all non-zero entries—meaning that the individual subjects are correlated (dependent)

More general situations...

- \Box General least squares represent a class of linear models for heteroscedastic data, however, with the known heteroscedastic structure—the matrix $\mathbb W$ is known from the experiment
- □ More general scenario involves situations where heteroscedastic data have some unknown variance structure (which needs to be estimated)
- □ Recall Assumption (A3) that specified the following conditions:

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

- □ The assumption above implies, that the matrix $\mathbb{W}^* = E[\sigma^2(\mathbf{X})\mathbf{X}\mathbf{X}^\top]$ is a real matrix with all elements being finite
- □ Thus, under the heteroscedastic model, we have $E[Y_i|X_i] = X_i^\top \beta$ and $Var[Y_i|X_i] = Var[\varepsilon_i|X_i] = \sigma^2(X_i)$

Consistency of the LSE estimates

The underlying model can be either assumed within the normal model framework or, alternatively, no normality is needed (some moment conditions are assumed instead)

□ Again, we are interested in the following parameters:

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

$$\begin{array}{c} \square \ \widehat{\beta}_n \longrightarrow \beta \text{ a.s. (in P), for } n \to \infty \\ \square \ \widehat{\theta}_n = \mathbf{I}^\top \widehat{\beta}_n \longrightarrow \theta \text{ a.s. (in P), for } n \to \infty \\ \square \ \widehat{\Theta}_n = \mathbb{L} \widehat{\beta}_n \longrightarrow \Theta, \text{ a.s. (in P), for } n \to \infty \end{array}$$

Assymptotic normality under heteroscedasticity

Under the assumptions stated in (A1), (A2), and (A3b) and, additionally, for $E[\varepsilon^2 X_j X_k] < \infty$ for j, k = 1, ..., p the following holds:

$$\begin{array}{c} \Box \quad \sqrt{n}(\widehat{\beta}_{n} - \beta) \stackrel{\mathcal{D}}{\longrightarrow} N_{p}(\beta, \sigma^{2} \mathbb{V} \mathbb{W}^{*} \mathbb{V}) \text{ for } n \to \infty \\ \Box \quad \sqrt{n}(\widehat{\theta}_{n} - \theta) \stackrel{\mathcal{D}}{\longrightarrow} N(0, \sigma^{2} I^{\top} \mathbb{V} \mathbb{W}^{*} \mathbb{V} I), \text{ as } n \to \infty \\ \Box \quad \sqrt{n}(\widehat{\Theta}_{n} - \Theta) \stackrel{\mathcal{D}}{\longrightarrow} N_{m}(\mathbf{0}, \sigma^{2} \mathbb{L} \mathbb{V} \mathbb{W}^{*} \mathbb{V} \mathbb{L}^{\top}), \text{ as } n \to \infty \\ \text{where } \mathbb{V} = \left[E(\mathbf{X} \mathbf{X}^{\top}) \right]^{-1} \text{ and } \mathbb{W}^{*} = E[\sigma^{2}(\mathbf{X}) \mathbf{X} \mathbf{X}^{\top}]$$

Note that $Var(\mathbf{X}\varepsilon) = E[\sigma^2(\mathbf{X})\mathbf{X}\mathbf{X}^{\top}]$ which equals to $\sigma^2 E[\mathbf{X}\mathbf{X}^{\top}] = \sigma^2 \mathbb{W}$ under homoscedasticity (A3a) and it equals to \mathbb{W}^* under heteroscedasticity (A3b)

Sandwich estimate of the variance

Consider the assumptions in (A1), (A2), and (A3b). Let, moreover, the following holds

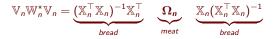
all for $j, k, s, l \in \{1, \dots, p\}$. Then the following holds:

 $n \mathbb{V}_n \mathbb{W}_n^* \mathbb{V}_n \xrightarrow{a.s.(P)} \mathbb{V} \mathbb{W}^* \mathbb{V}, \text{ for } n \to \infty$

where $\mathbb{W}_n^{\star} = \sum_{i=1}^n U_i^2 \mathbf{X}_i \mathbf{X}_i^{\top} = \mathbb{X}_n^{\top} \Omega_n \mathbb{X}_n$, where $U_i = Y_i - \widehat{Y}_i$ and $\Omega_n = diag(U_1^2, \dots, U_n^2)$

Sandwich estimate

□ the estimate for the variance covariance matrix VW^{*}V is the so-called **sandwich estimate** of the form



which is a (heteroscedastic) consistent estimate of the variance-covarance of the least squares estimate $\hat{\beta}_n$

- □ if we replace the matrix Ω_n with $\frac{n}{\nu_n}\Omega_n$ for some sequence $\{\nu_n\}_n$ such that $n/\nu_n \to 1$ as $n \to \infty$ the convergence still holds and ν_n is called the **degrees of freedom of the sandwich** estimate
- □ different options are used in the literature to define the sequence $\{\nu_n\}_n$ (White (1980); MacKinnon and White (1985); etc.)

Asymptotic inference under heteroscedasticity

□ for a consistent sandwich estimate $\mathbb{V}_n^{HC} = (\mathbb{X}_n^\top \mathbb{X}_n)^{-1} \mathbb{X}_n^\top \Omega_n \mathbb{X}_n (\mathbb{X}_n^\top \mathbb{X}_n)^{-1}$ of the covariance matrix of $\hat{\beta}_n$ we can define

$$T_n = \frac{I^\top \widehat{\beta}_n - I^\top \beta}{\sqrt{I^\top \mathbb{V}_n^{HC} I}}$$
$$Q_n = \frac{(\mathbb{L} \widehat{\beta}_n - \mathbb{L} \beta)^\top (\mathbb{L} \mathbb{V}_n^{HC} \mathbb{L}^\top)^{-1} (\mathbb{L} \widehat{\beta}_n - \mathbb{L} \beta)}{m}$$

- □ The statistic T_n follows (asymptotically) the normal distribution N(0,1)and the statistic mQ_n follows (again asymptotically) the χ^2 distribution with $m = rank(\mathbb{L})$ degrees of freedom (for $n \to \infty$)
- □ Note that the results are analogous to those obtained for the homoscedastic situation where MSe(X^TX)⁻¹ is replaced by the sandwich estimate V^{HC}_n
- □ the statistics T_n and Q_n can be directly used to perform statistical inference—i.e., to constract a confidence interval/region or to test some set of hypotheses

Summary

Linear regression models

- Normal linear model with homoscedastic errors
- □ Linear model without normality assumptions (A3a/A3b)
- General linear model (with and without the normality assumption)

Consistent LSE/MLE estimates

- consistent estimates of the mean and variance parameters
- □ the mean parameter estimates are normally distributed (normal model)
- the mean estimates are asymptotically normal (model without normality)
- consistent estimates of the variance parameter/parameters

Statistical inference

- primarily about the mean parameters and their linear combinations
- exact and approximate (asymptotic) confidence intervals (regions)
- □ statistical tests (hull and alternative hypotheses)