

FACULTY
OF MATHEMATICS
AND PHYSICS
Charles University

MATHEMATICAL STATISTICS 3

NMST 424

Course notes

Last updated: April 30, 2024

I would like to thank Stanislav Nagy for finding many typos and misprints. I would like to thank all the students and colleagues who have helped me in improving the text.

Contents

1	Clippings from the asymptotic theory	1
1.1	The convergence of random vectors	1
1.2	Δ -theorem	5
1.3	Moment estimators	9
1.4	Confidence intervals and asymptotic variance-stabilising transformation	13
2	Maximum likelihood methods	16
2.1	Asymptotic normality of maximum likelihood estimator	16
2.2	Asymptotic efficiency of maximum likelihood estimators	21
2.3	Estimation of the asymptotic variance matrix	22
2.4	Asymptotic tests (without nuisance parameters)	23
2.5	Asymptotic confidence sets	25
2.6	Asymptotic tests with nuisance parameters	26
2.7	Profile likelihood	33
2.8	Some notes on maximum likelihood in case of not i.i.d. random vectors	37
2.9	Conditional and marginal likelihood	40
3	M- and Z-estimators	46
3.1	Identifiability of parameters via M - and/or Z -estimators	48
3.2	Asymptotic distribution of Z -estimators	49
3.3	Likelihood under model misspecification	54
3.4	Asymptotic normality of M -estimators defined by convex minimization	56
4	M-estimators and Z-estimators in robust statistics	59
4.1	Robust estimation of location	60
4.2	Robust studentized M/Z -estimators of location	62
4.3	Robust estimation in linear models	64
4.3.1	The least squares method	64
4.3.2	Method of the least absolute deviation	64
4.3.3	Huber estimator of regression	66
4.3.4	Studentized Huber estimator of regression	67
5	Quantile regression	68
5.1	Identification of quantiles	68
5.2	Regression quantiles	70
5.3	Interpretation of the regression quantiles	74

5.4	Inference for regression quantiles	75
5.5	Asymptotic normality of sample quantiles	76
6	EM-algorithm	76
6.1	General description of the EM-algorithm	79
6.2	Convergence of the EM-algorithm	81
6.3	Rate of convergence of EM-algorithm	82
6.4	The EM algorithm in exponential families	83
7	Missing data	84
7.1	Basic concepts for the mechanism of missing	85
7.2	Methods for dealing with missing data	87

Last update: April 30, 2024

1 Clippings from the asymptotic theory

1.1 The convergence of random vectors

Let \mathbf{X} be a k -dimensional random vector (with the cumulative distribution function $F_{\mathbf{X}}$) and $\{\mathbf{X}_n\}_{n=1}^{\infty}$ be a sequence of k -dimensional random vectors (with the cumulative distribution functions $F_{\mathbf{X}_n}$).

Definition. We say that $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$ (i.e. \mathbf{X}_n converges *in distribution* to \mathbf{X}), if

$$\lim_{n \rightarrow \infty} F_{\mathbf{X}_n}(\mathbf{x}) = F_{\mathbf{X}}(\mathbf{x})$$

for each point \mathbf{x} of the continuity of $F_{\mathbf{X}}$.

Example 1. Let U be a random variable with a uniform distribution on the interval $(0, 1)$. Put $X_n = U/n$. Show that $X_n \xrightarrow[n \rightarrow \infty]{d} 0$. But at the same time $F_{X_n}(0)$ does not converge to $F_X(0)$, where F_X is the cumulative distribution corresponding to the random variable that is equal to zero almost surely.

Let d be a metric in \mathbb{R}^k , e.g. the Euclidean metric $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{j=1}^k (x_j - y_j)^2}$.

Definition. We say that

- $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{X}$ (i.e. \mathbf{X}_n converges *in probability* to \mathbf{X}), if

$$\forall \varepsilon > 0 \lim_{n \rightarrow \infty} \mathbb{P} \left[\omega : d(\mathbf{X}_n(\omega), \mathbf{X}(\omega)) > \varepsilon \right] = 0;$$

- $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mathbf{X}$ (i.e. \mathbf{X}_n converges *almost surely* to \mathbf{X}), if

$$\mathbb{P} \left[\omega : \lim_{n \rightarrow \infty} d(\mathbf{X}_n(\omega), \mathbf{X}(\omega)) = 0 \right] = 1.$$

Remark 1. For random vectors the convergence in probability and almost surely can be defined also component-wise. That is let $\mathbf{X}_n = (X_{n1}, \dots, X_{nk})^{\top}$ and $\mathbf{X} = (X_1, \dots, X_k)^{\top}$. Then

$$\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{X} \quad (\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mathbf{X}) \quad \text{if} \quad X_{nj} \xrightarrow[n \rightarrow \infty]{P} X_j \quad (X_{nj} \xrightarrow[n \rightarrow \infty]{\text{a.s.}} X_j), \quad \forall j = 1, \dots, k.$$

But this is not true for the convergence in distribution for which we have the Cramér-Wold theorem that states

$$\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X} \iff \boldsymbol{\lambda}^{\top} \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \boldsymbol{\lambda}^{\top} \mathbf{X}, \quad \forall \boldsymbol{\lambda} \in \mathbb{R}^k.$$

Theorem 1. (Continuous Mapping Theorem, CMT) Let $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^m$ be continuous in each point of an open set $C \subset \mathbb{R}^k$ such that $\mathbb{P}(\mathbf{X} \in C) = 1$. Then

$$(i) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{a.s.} \mathbf{X} \Rightarrow \mathbf{g}(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{a.s.} \mathbf{g}(\mathbf{X});$$

$$(ii) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{X} \Rightarrow \mathbf{g}(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{P} \mathbf{g}(\mathbf{X});$$

$$(iii) \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X} \Rightarrow \mathbf{g}(\mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{d} \mathbf{g}(\mathbf{X}).$$

Proof. (i) *Almost sure convergence.*

$$\begin{aligned} & \mathbb{P} \left[\omega : \lim_{n \rightarrow \infty} d(\mathbf{g}(\mathbf{X}_n(\omega)), \mathbf{g}(\mathbf{X}(\omega))) = 0 \right] \\ & \geq \mathbb{P} \left[\omega : \lim_{n \rightarrow \infty} d(\mathbf{g}(\mathbf{X}_n(\omega)), \mathbf{g}(\mathbf{X}(\omega))) = 0, \mathbf{X}(\omega) \in C \right] \\ & \geq \mathbb{P} \left[\omega : \lim_{n \rightarrow \infty} d(\mathbf{X}_n(\omega), \mathbf{X}(\omega)) = 0, \mathbf{X}(\omega) \in C \right] = 1, \end{aligned}$$

as C is an open set and $\mathbb{P}(\mathbf{X} \in C) = 1$.

(ii) *Convergence in probability.* Let $\varepsilon > 0$. Then for each $\delta > 0$

$$\begin{aligned} & \mathbb{P} \left[\omega : d(\mathbf{g}(\mathbf{X}_n(\omega)), \mathbf{g}(\mathbf{X}(\omega))) > \varepsilon \right] \\ & \leq \mathbb{P} \left[d(\mathbf{g}(\mathbf{X}_n), \mathbf{g}(\mathbf{X})) > \varepsilon, d(\mathbf{X}_n, \mathbf{X}) \leq \delta \right] + \mathbb{P} \left[d(\mathbf{X}_n, \mathbf{X}) > \delta \right] \\ & \leq \mathbb{P} \left[\mathbf{X} \in B^\delta \right] + \underbrace{\mathbb{P} \left[d(\mathbf{X}_n, \mathbf{X}) > \delta \right]}_{\rightarrow 0, \forall \delta > 0}, \end{aligned}$$

where $B^\delta = \{ \mathbf{x} \in \mathbb{R}^k; \exists \mathbf{y} \in \mathbb{R}^k : d(\mathbf{x}, \mathbf{y}) \leq \delta, d(\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{y})) > \varepsilon \}$. Further

$$\begin{aligned} \mathbb{P} \left[\mathbf{X} \in B^\delta \right] &= \mathbb{P} \left[\mathbf{X} \in B^\delta, \mathbf{X} \in C \right] + \mathbb{P} \left[\mathbf{X} \in B^\delta, \mathbf{X} \notin C \right] \\ &= \mathbb{P} \left[\mathbf{X} \in B^\delta \cap C \right] + 0 \end{aligned}$$

and $\mathbb{P} \left[\mathbf{X} \in B^\delta \cap C \right]$ can be made arbitrarily small as $B^\delta \cap C \rightarrow \emptyset$ for $\delta \searrow 0$.

(iii) See for instance the proof of Theorem 13.6 in [Lachout \[2004\]](#). □

Theorem 2. (Cramér-Slutsky, CS) Let $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$, $\mathbf{Y}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{c}$, then

$$(i) \mathbf{X}_n + \mathbf{Y}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X} + \mathbf{c};$$

$$(ii) \mathbf{Y}_n \mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{c} \mathbf{X},$$

where \mathbf{Y}_n can be a sequence of random variables or vectors or matrices of appropriate dimensions (\mathbb{R} or \mathbb{R}^k or $\mathbb{R}^{m \times k}$) and analogously \mathbf{c} can be either a number or a vector or a matrix of an appropriate dimension.

Proof. Note that it is sufficient to prove

$$(\mathbf{X}_n, \mathbf{Y}_n) \xrightarrow[n \rightarrow \infty]{d} (\mathbf{X}, \mathbf{c}). \quad (1)$$

Then the statement of the theorem follows from Continuous Mapping Theorem (Theorem 1).

To prove (1) note that

$$d((\mathbf{X}_n, \mathbf{Y}_n), (\mathbf{X}_n, \mathbf{c})) = d(\mathbf{Y}_n, \mathbf{c}) \xrightarrow[n \rightarrow \infty]{P} 0.$$

Thus by Theorem 13.7 in Lachout [2004] or Theorem 2.7 (iv) of van der Vaart [2000] it is sufficient to show that $(\mathbf{X}_n, \mathbf{c}) \xrightarrow[n \rightarrow \infty]{d} (\mathbf{X}, \mathbf{c})$. But this follows immediately with the help of the Cramér-Wold theorem. \square

Definition 1. Let $\{\mathbf{X}_n\}_{n=1}^{\infty}$ be a sequence of random vectors and $\{r_n\}_{n=1}^{\infty}$ a sequence of positive constants. We write that

- (i) $\mathbf{X}_n = o_P\left(\frac{1}{r_n}\right)$, if $(r_n \mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{P} \mathbf{0}_k$, where $\mathbf{0}_k = (0, \dots, 0)^\top$ is a zero point in \mathbb{R}^k ;
- (ii) $\mathbf{X}_n = O_P\left(\frac{1}{r_n}\right)$, if

$$\forall \varepsilon > 0 \exists K < \infty \exists n_0 \in \mathbb{N} \sup_{n \geq n_0} \mathbb{P}\left(r_n \|\mathbf{X}_n\| > K\right) < \varepsilon,$$

where $\|\cdot\|$ stands for instance for the Euclidean norm.

When $\mathbf{X}_n = O_P(1)$ then some authors say that $\{\mathbf{X}_n\}$ is (asymptotically) *bounded* in probability*. When $\mathbf{X}_n = o_P(1)$ then it is often said that $\{\mathbf{X}_n\}$ is (asymptotically) *negligible* in probability.

Remark 2. Note that

- (i) $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$ implies $\mathbf{X}_n = O_P(1)$ (Prohorov's theorem, Portmanteau theorem, see e.g. Chapters 2.1 van der Vaart [2000]);
- (ii) $\mathbf{X}_n \xrightarrow[n \rightarrow \infty]{P} \mathbf{0}$ implies $\mathbf{X}_n = o_P(1)$;
- (iii) $(r_n \mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{P} \mathbf{X}$ or $(r_n \mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{d} \mathbf{X}$ implies $\mathbf{X}_n = O_P\left(\frac{1}{r_n}\right)$.
- (iv) If $r_n \rightarrow \infty$ and $\mathbf{X}_n = O_P\left(\frac{1}{r_n}\right)$, then $\mathbf{X}_n = o_P(1)$.

* *omezená v pravděpodobnosti*

Proof of (iv). Note that it is sufficient to prove that for each $\varepsilon > 0$ and each $\eta > 0$ for all sufficiently large n it holds that $\mathbf{P}(\|\mathbf{X}_n\| > \varepsilon) < \eta$.

Note that $\mathbf{X}_n = O_P(\frac{1}{r_n})$ implies there exists a finite constant K and $n_0 \in \mathbb{N}$ such that

$$\sup_{n \geq n_0} \mathbf{P}(r_n \|\mathbf{X}_n\| > K) < \varepsilon.$$

The statement now follows from the fact that

$$\mathbf{P}(\|\mathbf{X}_n\| > \varepsilon) = \mathbf{P}(r_n \|\mathbf{X}_n\| > \varepsilon r_n) < \eta$$

for all n such that $\varepsilon r_n > K$. □

Suppose that $\mathbf{X}_1, \mathbf{X}_2, \dots$ are independent and identically distributed random vectors with a finite variance matrix. Then the law of large numbers implies

$$\bar{\mathbf{X}}_n = \mathbf{E} \mathbf{X}_1 + o_P(1).$$

With the help of the central limit theorem one can be even more specific about the remainder term and show that

$$\bar{\mathbf{X}}_n = \mathbf{E} \mathbf{X}_1 + O_P\left(\frac{1}{\sqrt{n}}\right).$$

Remark 3. Further note that the calculus with the random quantities $o_P(1)$ and $O_P(1)$ is analogous to the calculus with the (deterministic) quantities $o(1)$ and $O(1)$ in mathematical analysis. Thus, among others it holds that

(i) $o_P(1) + o_P(1) = o_P(1)$;

(ii) $o_P(1) O_P(1) = o_P(1)$;

(iii) $O_P(1) O_P(1) = O_P(1)$;

(iv) $o_P(1) + O_P(1) = O_P(1)$;

Proof of (ii). Let $\{\mathbf{X}_n\}, \{\mathbf{Y}_n\}$ be such that $\mathbf{X}_n = O_P(1), \mathbf{Y}_n = o_P(1)$ and $\mathbf{Y}_n \mathbf{X}_n$ makes sense. Let $\varepsilon > 0$ be given and consider for instance the Euclidean norm (for other norms the proof would go through up to a multiplicative constant in some of the arguments). Then one can find $K < \infty$ and $n_0 \in \mathbb{N}$ such that $\sup_{n \in \mathbb{N}_0} \mathbf{P}(\|\mathbf{X}_n\| > K) < \frac{\varepsilon}{2}$. Thus for all sufficiently large $n \in \mathbb{N}$

$$\begin{aligned} \mathbf{P}(\|\mathbf{Y}_n \mathbf{X}_n\| > \varepsilon) &\leq \mathbf{P}(\|\mathbf{Y}_n \mathbf{X}_n\| > \varepsilon, \|\mathbf{X}_n\| \leq K) + \mathbf{P}(\|\mathbf{X}_n\| > K) \\ &\leq \mathbf{P}(\|\mathbf{Y}_n\| > \frac{\varepsilon}{K}) + \frac{\varepsilon}{2} \leq \varepsilon, \end{aligned}$$

as $\mathbf{Y}_n = o_P(1)$.

We recommend the reader to prove the remaining statements as an exercise. □

For more details about the calculus with $o_P(1)$ and $O_P(1)$ see for instance Chapter 3.4 of Jiang [2010].

1.2 Δ -theorem

Let $\mathbf{T}_n = (T_{n1}, \dots, T_{np})^\top$ be a p -dimensional random vector that converges to the constant $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top$ and $\mathbf{g} = (g_1, \dots, g_m)^\top$ be a function from (a subset of) \mathbb{R}^p to \mathbb{R}^m . Denote the Jacobi matrix of the function \mathbf{g} at the point \mathbf{x} as $\mathbb{D}_{\mathbf{g}}(\mathbf{x})$, i.e.

$$\mathbb{D}_{\mathbf{g}}(\mathbf{x}) = \begin{pmatrix} \nabla g_1(\mathbf{x}) \\ \vdots \\ \nabla g_m(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \frac{\partial g_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial g_1(\mathbf{x})}{\partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_m(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial g_m(\mathbf{x})}{\partial x_p} \end{pmatrix}.$$

Theorem 3. (Δ -theorem) Let $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = O_P(1)$. Further $\mathbf{g} : A \rightarrow \mathbb{R}^m$, where $A \subset \mathbb{R}^p$, $\boldsymbol{\mu}$ is an interior point of A and the first-order partial derivatives of \mathbf{g} are continuous in a neighbourhood of $\boldsymbol{\mu}$. Then

(i) $\sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\mu})) - \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = o_P(1)^*$;

(ii) moreover if $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma})$, then

$$\sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\mu})) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_m(\mathbf{0}_m, \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu}) \boldsymbol{\Sigma} \mathbb{D}_{\mathbf{g}}^\top(\boldsymbol{\mu})). \quad (2)$$

Proof. Statement (i): For $j \in \{1, \dots, m\}$ consider $g_j : A \rightarrow \mathbb{R}$ (the j -th coordinate of the function \mathbf{g}). From the assumptions of the theorem there exists a neighbourhood $\mathcal{U}_\delta(\boldsymbol{\mu})$ of the point $\boldsymbol{\mu}$ such that the function g_j has continuous partial derivatives in this neighbourhood. Further $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = O_P(1)$ implies $\mathbf{T}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$ (see for instance Remark 2(iv)), which yields that $\mathbb{P}(\mathbf{T}_n \in \mathcal{U}_\delta(\boldsymbol{\mu})) \xrightarrow[n \rightarrow \infty]{} 1$. Thus without loss of generality one can assume that $\mathbf{T}_n \in \mathcal{U}_\delta(\boldsymbol{\mu})$. Using this together with the mean value theorem there exists $\boldsymbol{\mu}_n^{j*}$ which lies between \mathbf{T}_n and $\boldsymbol{\mu}$ such that

$$\begin{aligned} \sqrt{n}(g_j(\mathbf{T}_n) - g_j(\boldsymbol{\mu})) &= \nabla g_j(\boldsymbol{\mu}_n^{j*})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) \\ &= \nabla g_j(\boldsymbol{\mu})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) + [\nabla g_j(\boldsymbol{\mu}_n^{j*}) - \nabla g_j(\boldsymbol{\mu})]\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}). \end{aligned} \quad (3)$$

Further $\mathbf{T}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$ implies that $\boldsymbol{\mu}_n^{j*} \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$. Now the continuity of the partial derivatives of g_j in $\mathcal{U}_\delta(\boldsymbol{\mu})$ and CMT (Theorem 1) imply that

$$\nabla g_j(\boldsymbol{\mu}_n^{j*}) - \nabla g_j(\boldsymbol{\mu}) = o_P(1),$$

* $\mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu})$ is sometimes called also the asymptotic linear approximation of $\sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\mu}))$.

which together with $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = O_P(1)$ gives

$$[\nabla g_j(\boldsymbol{\mu}_n^{j*}) - \nabla g_j(\boldsymbol{\mu})] \sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = o_P(1). \quad (4)$$

Now combining (3) and (4) yields that for each $j = 1, \dots, m$

$$\sqrt{n}(g_j(\mathbf{T}_n) - g_j(\boldsymbol{\mu})) = \nabla g_j(\boldsymbol{\mu})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) + o_P(1),$$

which implies the first statement of the theorem.

Statement (i) - the proof done at the lecture: For $j \in \{1, \dots, m\}$ consider $g_j : A \rightarrow \mathbb{R}$ (the j -th coordinate of the function \mathbf{g}). Define the following function $h_j : A \rightarrow \mathbb{R}$ as

$$h_j(\mathbf{x}) = \begin{cases} \frac{g_j(\mathbf{x}) - g_j(\boldsymbol{\mu}) - \nabla g_j(\boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})}{\|\mathbf{x} - \boldsymbol{\mu}\|}, & \mathbf{x} \neq \boldsymbol{\mu}, \\ 0, & \mathbf{x} = \boldsymbol{\mu}. \end{cases}$$

Note that h_j is continuous in $\boldsymbol{\mu}$.

Further $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = O_P(1)$ implies $\mathbf{T}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\mu}$ (see for instance Remark 2(iv)). Now using CMT (Theorem 1) implies $h_j(\mathbf{T}_n) = o_P(1)$. Thus

$$g_j(\mathbf{T}_n) - g_j(\boldsymbol{\mu}) - \nabla g_j(\boldsymbol{\mu})(\mathbf{T}_n - \boldsymbol{\mu}) = o_P(1) \|\mathbf{T}_n - \boldsymbol{\mu}\|,$$

which together with $\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = O_P(1)$ and Remark 3(ii) gives

$$\sqrt{n}(g_j(\mathbf{T}_n) - g_j(\boldsymbol{\mu})) - \nabla g_j(\boldsymbol{\mu})(\mathbf{T}_n - \boldsymbol{\mu}) = o_P(1) \|\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu})\| = o_P(1).$$

Thus one can conclude that

$$\sqrt{n}(g_j(\mathbf{T}_n) - g_j(\boldsymbol{\mu})) = \nabla g_j(\boldsymbol{\mu})\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) + o_P(1),$$

which implies the first statement of the theorem.

Statement (ii): By the first statement of the theorem one gets

$$\sqrt{n}(\mathbf{g}(\bar{\mathbf{X}}_n) - \mathbf{g}(\boldsymbol{\mu})) = \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) + o_P(1)$$

Now for the term $\mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu})$ one can use the second statement of CS (Theorem 2) with $\mathbf{Y}_n = \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})$ and $\mathbf{X}_n = \sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu})$. Further, using now the first statement of CS with $\mathbf{c} = \mathbf{0}_m$ one can see that adding the term $o_P(1)$ does not alter the asymptotic distribution of $\mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu})$. \square

In the most common applications of Δ -theorem one often takes $\mathbf{T}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, where $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent and identically distributed. Then $\boldsymbol{\mu} = \mathbb{E} \mathbf{X}_1$ and the standard central limit theorem gives the asymptotic normality

$$\sqrt{n}(\mathbf{T}_n - \boldsymbol{\mu}) = \sqrt{n}(\bar{\mathbf{X}}_n - \mathbb{E} \mathbf{X}_1) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma}),$$

where $\Sigma = \text{var}(\mathbf{X}_i)$.

Note that then Theorem 3(i) implies that

$$\sqrt{n}(\mathbf{g}(\bar{\mathbf{X}}_n) - \mathbf{g}(\boldsymbol{\mu})) = \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu}) \sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) + o_P(1) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{Z}_i + o_P(1),$$

where

$$\mathbf{Z}_i = \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu}), \quad i \in \{1, \dots, n\}$$

are independently distributed random vectors. Note that then the central limit theorem together with the Cramér-Slutsky theorem (Theorem 2(i)) implies that

$$\sqrt{n}(\mathbf{g}(\bar{\mathbf{X}}_n) - \mathbf{g}(\boldsymbol{\mu})) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(\mathbf{0}, \text{var}(\mathbf{Z}_1)).$$

Thus the asymptotic variance

$$\text{avar}(\mathbf{g}(\bar{\mathbf{X}}_n)) = \frac{1}{n} \text{var}(\mathbf{Z}_1)$$

can be easily estimated as

$$\widehat{\text{avar}}(\mathbf{g}(\bar{\mathbf{X}}_n)) = \frac{1}{n} \mathbb{S}_Z^2,$$

where

$$\mathbb{S}_Z^2 = \frac{1}{n-1} \sum_{i=1}^n \hat{\mathbf{Z}}_i \hat{\mathbf{Z}}_i^\top$$

is the sample variance matrix of the ‘estimated’ \mathbf{Z}_i given by

$$\hat{\mathbf{Z}}_i = \mathbb{D}_{\mathbf{g}}(\bar{\mathbf{X}}_n)(\mathbf{X}_i - \bar{\mathbf{X}}_n),$$

as $\sum_{i=1}^n \hat{\mathbf{Z}}_i = \mathbf{0}$.

Remark 4. Instead of the continuity of the partial derivatives in a neighbourhood of $\boldsymbol{\mu}$, it would be sufficient to assume the existence of the total differential of the function \mathbf{g} at the point $\boldsymbol{\mu}$ (see the alternative proof of (i) done at the lecture).

Sometimes instead of (2) we write shortly $\mathbf{g}(\mathbf{T}_n) \stackrel{\text{as}}{\approx} \mathbf{N}_m(\mathbf{g}(\boldsymbol{\mu}), \frac{1}{n} \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu}) \Sigma \mathbb{D}_{\mathbf{g}}^\top(\boldsymbol{\mu}))$. The quantity $\frac{1}{n} \mathbb{D}_{\mathbf{g}}(\boldsymbol{\mu}) \Sigma \mathbb{D}_{\mathbf{g}}^\top(\boldsymbol{\mu})$ is then called the **asymptotic variance** matrix of $\mathbf{g}(\mathbf{T}_n)$ and it is denoted as $\text{avar}(\mathbf{g}(\mathbf{T}_n))$. Note that the asymptotic variance has to be understood as the **variance of the asymptotic distribution**, but not as a limiting variance.

As the following three examples show for a sequence of random variables $\{Y_n\}$ the asymptotic variance $\text{avar}(Y_n)$ may exist even if $\text{var}(Y_n)$ does not exist for any $n \in \mathbb{N}$. Further even if $\text{var}(Y_n)$ exists, then it **does not hold that** $\text{var}(Y_n)/\text{avar}(Y_n) \rightarrow 1$ as $n \rightarrow \infty$.

Example 2. Let $X \sim \mathbf{N}(0, 1)$ and $\{\varepsilon_n\}$ be a sequence of random variables independent with X such that

$$\mathbf{P}(\varepsilon_n = -\sqrt{n}) = \frac{1}{2n}, \quad \mathbf{P}(\varepsilon_n = 0) = 1 - \frac{1}{n}, \quad \mathbf{P}(\varepsilon_n = \sqrt{n}) = \frac{1}{2n}.$$

Define $Y_n = X + \varepsilon_n$ and show that $Y_n \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, 1)$. Thus $\text{avar}(Y_n) = 1$. On the other hand $\text{var}(Y_n) = 2$ for each $n \in \mathbb{N}$.

Example 3. A random sample X_1, \dots, X_n from a zero-mean distribution with finite and positive variance. Find the asymptotic distribution of $Y_n = \bar{X}_n \exp\{\bar{X}_n^3\}$. Further compare $\text{var}(Y_n)$ and $\text{avar}(Y_n)$ when X_1 is distributed as $\mathbf{N}(0, 1)$.

Example 4. Suppose you have a random sample X_1, \dots, X_n from a Bernoulli distribution with parameter p_X and you are interested in estimating the logarithm of the odd, i.e. $\theta_X = \log\left(\frac{p_X}{1-p_X}\right)$. Compare the variance and the asymptotic variance of $\hat{\theta}_X = \log\left(\frac{\bar{X}_n}{1-\bar{X}_n}\right)$.

Example 5. Suppose you have two independent random samples from Bernoulli distribution. Derive the asymptotic distribution of the logarithm of odds-ratio.

Example 6. Suppose we observe independent identically distributed random vectors

$$\begin{pmatrix} X_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ Y_n \end{pmatrix}$$

and denote $\rho = \frac{\text{cov}(X_1, Y_1)}{\sqrt{\text{var}(X_1)\text{var}(Y_1)}}$ the (Pearson's) correlation coefficient. Consider the sample correlation coefficient given by

$$\hat{\rho}_n = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2 \sum_{i=1}^n (Y_i - \bar{Y}_n)^2}}.$$

With the help of Theorem 3(i) derive (the asymptotic representation)

$$\sqrt{n}(\hat{\rho}_n - \rho) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\tilde{X}_i \tilde{Y}_i - \frac{\rho}{2} \tilde{X}_i^2 - \frac{\rho}{2} \tilde{Y}_i^2] + o_P(1),$$

where $\tilde{X}_i = \frac{X_i - \mathbf{E} X_1}{\sqrt{\text{var}(X_1)}}$ and $\tilde{Y}_i = \frac{Y_i - \mathbf{E} Y_1}{\sqrt{\text{var}(Y_1)}}$ are standardized versions of X_i and Y_i . Conclude that

$$\sqrt{n}(\hat{\rho}_n - \rho) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, \text{var}(Z_i)),$$

where $Z_i = \tilde{X}_i \tilde{Y}_i - \frac{\rho}{2} \tilde{X}_i^2 - \frac{\rho}{2} \tilde{Y}_i^2$. Derive the asymptotic distribution under the independence of X_i and Y_i and suggest a test of independence.

Further show that if $(X_i, Y_i)^\top$ follows the bivariate normal distribution then

$$\sqrt{n}(\hat{\rho}_n - \rho) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, (1 - \rho^2)^2).$$

Find the (asymptotic) variance stabilising transformation for $\hat{\rho}_n$ (see Chapter 1.4) and derive the confidence interval for ρ .

Example 7. Consider a random sample from the Bernoulli distribution with the parameter p_X . Derive the asymptotic distribution of the estimator of $\theta_X = p_X(1 - p_X)$ (variance of the Bernoulli distribution) given by $\hat{\theta}_n = \frac{n}{n-1} \bar{X}_n(1 - \bar{X}_n)$.

Example 8. Suppose that we observe X_1, \dots, X_n of a moving average sequence of order 1 given by

$$X_t = Y_t + \theta Y_{t-1}, \quad t \in \mathbb{Z},$$

where $\{Y_t, t \in \mathbb{Z}\}$ is a white noise sequence such that $\mathbb{E} Y_t = 0$ and $\text{var}(Y_t) = \sigma^2$.

Using the fact that the autocorrelation function at lag 1 satisfies

$$r(1) = \frac{\theta}{1 + \theta^2}$$

derive the estimator of θ and find its asymptotic distribution.

Hint. Note that by Bartlett's formula

$$\sqrt{n} (\hat{r}_n(1) - r(1)) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, \sigma^2(\theta)),$$

where

$$\sigma^2(\theta) = 1 - 3\left(\frac{\theta}{1+\theta^2}\right)^2 + 4\left(\frac{\theta}{1+\theta^2}\right)^4.$$

The end of
class 2
(23. 2. 2024)

1.3 Moment estimators

Suppose that the random vector \mathbf{X} has a density $f(\mathbf{x}; \boldsymbol{\theta})$ with respect to a σ -finite measure μ and that the density is known up to unknown p -dimensional parameter $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^\top \in \Theta$. Let $\boldsymbol{\theta}_X$ be the true value* of this unknown parameter. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from this distribution and t_1, \dots, t_p be given real functions. For instance if the observations are one-dimensional one can take $t_j(x) = x^j$, $j \in \{1, \dots, p\}$. For $j \in \{1, \dots, p\}$ define the function $\tau_j : \Theta \rightarrow \mathbb{R}$ as

$$\tau_j(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} t_j(\mathbf{X}_1) = \int t_j(\mathbf{x}) f(\mathbf{x}; \boldsymbol{\theta}) d\mu(\mathbf{x}), \quad j \in \{1, \dots, p\}.$$

Then the moment estimator[†] $\hat{\boldsymbol{\theta}}_n$ of the parameter $\boldsymbol{\theta}$ is a solution to the estimating equations

$$\frac{1}{n} \sum_{i=1}^n t_1(\mathbf{X}_i) = \tau_1(\hat{\boldsymbol{\theta}}_n), \dots, \frac{1}{n} \sum_{i=1}^n t_p(\mathbf{X}_i) = \tau_p(\hat{\boldsymbol{\theta}}_n).$$

Example 9. Let X_1, \dots, X_n be a random sample from the Beta distribution with the density $f(x; \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \mathbb{1}\{x \in (0, 1)\}$. Consider $t_1(x) = x$ and $t_2(x) = x^2$. Then

$$\mathbb{E} X_1 = \frac{\alpha}{\alpha + \beta}, \quad \mathbb{E} X_1^2 = \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)}.$$

* skutečná hodnota † Momentový odhad

Thus the estimating equations are

$$\frac{1}{n} \sum_{i=1}^n X_i = \frac{\hat{\alpha}}{\hat{\alpha} + \hat{\beta}}, \quad \frac{1}{n} \sum_{i=1}^n X_i^2 = \frac{\hat{\alpha}(\hat{\alpha} + 1)}{(\hat{\alpha} + \hat{\beta})(\hat{\alpha} + \hat{\beta} + 1)}.$$

Now denote $\boldsymbol{\tau}(\alpha, \beta) = (\tau_1(\alpha, \beta), \tau_2(\alpha, \beta))^\top$, where

$$\tau_1(\alpha, \beta) = \frac{\alpha}{\alpha + \beta}, \quad \tau_2(\alpha, \beta) = \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)}.$$

Thus one can rewrite the estimating equations as

$$\boldsymbol{\tau}(\hat{\alpha}, \hat{\beta}) = \left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2 \right)^\top.$$

Now provided that the inverse function $\boldsymbol{\tau}^{-1}$ exists one can write

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \boldsymbol{\tau}^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2 \right)^\top.$$

and use Δ -theorem to derive the asymptotic distribution of the estimator $\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$.

Now in the general situation put

$$\mathbf{T}_n = \left(\frac{1}{n} \sum_{i=1}^n t_1(\mathbf{X}_i), \dots, \frac{1}{n} \sum_{i=1}^n t_p(\mathbf{X}_i) \right)^\top \quad (5)$$

and define the mapping $\boldsymbol{\tau} : \Theta \mapsto \mathbb{R}^p$ as $\boldsymbol{\tau}(\boldsymbol{\theta}) = (\tau_1(\boldsymbol{\theta}), \dots, \tau_p(\boldsymbol{\theta}))^\top$. Note that provided there exists an inverse mapping $\boldsymbol{\tau}^{-1}$ one can write

$$\sqrt{n} \left(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X \right) = \sqrt{n} \left(\boldsymbol{\tau}^{-1}(\mathbf{T}_n) - \boldsymbol{\tau}^{-1}(\boldsymbol{\tau}(\boldsymbol{\theta}_X)) \right). \quad (6)$$

Thus the asymptotic normality of the moment estimator $\hat{\boldsymbol{\theta}}_n$ would follow by the Δ -theorem (Theorem 3) with $\mathbf{g} = \boldsymbol{\tau}^{-1}$. This is formalized in the following theorem.

Theorem 4. *Let $\boldsymbol{\theta}_X$ be an interior point of Θ and $\max_{j \in \{1, \dots, p\}} \text{var}_{\boldsymbol{\theta}_X}(t_j(\mathbf{X}_1)) < \infty$. Further let the function $\boldsymbol{\tau}$ be one-to-one and have continuous first-order partial derivatives in a neighbourhood of $\boldsymbol{\theta}_X$. Finally let the Jacobi matrix $\mathbb{D}_{\boldsymbol{\tau}}(\boldsymbol{\theta}_X)$ be regular. Then*

$$\sqrt{n} \left(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X \right) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{D}_{\boldsymbol{\tau}^{-1}}(\boldsymbol{\theta}_X) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{D}_{\boldsymbol{\tau}^{-1}}(\boldsymbol{\theta}_X)]^\top),$$

where $\boldsymbol{\Sigma}(\boldsymbol{\theta}_X) = \text{var}_{\boldsymbol{\theta}_X}(t_1(\mathbf{X}_1), \dots, t_p(\mathbf{X}_1))$.

Proof. By the assumptions of the theorem and the inverse function theorem (Theorem A13) there exists an open neighbourhood U containing $\boldsymbol{\theta}_X$ and an open neighbourhood V containing $\boldsymbol{\tau}(\boldsymbol{\theta}_X)$ such that $\boldsymbol{\tau} : U \rightarrow V$ is a differentiable bijection with a differentiable inverse $\boldsymbol{\tau}^{-1} : V \rightarrow U$. Further note that \mathbf{T}_n defined in (5) satisfies $\mathbb{P}(\mathbf{T}_n \in V) \xrightarrow[n \rightarrow \infty]{} 1$. Thus one can use (6) and apply the Δ -theorem (Theorem 3) with $\mathbf{g} = \boldsymbol{\tau}^{-1}$, $\boldsymbol{\mu} = \boldsymbol{\tau}(\boldsymbol{\theta}_X)$ and $A = V$ to get

$$\sqrt{n} \left(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X \right) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p \left(\mathbf{0}, \mathbb{D}_{\boldsymbol{\tau}^{-1}}(\boldsymbol{\tau}(\boldsymbol{\theta}_X)) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{D}_{\boldsymbol{\tau}^{-1}}(\boldsymbol{\tau}(\boldsymbol{\theta}_X))]^\top \right).$$

The statement of the theorem now follows from the identity

$$\mathbb{D}_{\boldsymbol{\tau}^{-1}}(\boldsymbol{\tau}(\boldsymbol{\theta}_X)) = \mathbb{D}_{\boldsymbol{\tau}}^{-1}(\boldsymbol{\theta}_X).$$

□

The asymptotic variance of $\widehat{\boldsymbol{\theta}}_n$ is usually estimated as

$$\frac{1}{n} \mathbb{D}_{\boldsymbol{\tau}}^{-1}(\widehat{\boldsymbol{\theta}}_n) \widehat{\boldsymbol{\Sigma}}_n [\mathbb{D}_{\boldsymbol{\tau}}^{-1}(\widehat{\boldsymbol{\theta}}_n)]^\top,$$

where as $\widehat{\boldsymbol{\Sigma}}_n$ one can take either $\boldsymbol{\Sigma}(\widehat{\boldsymbol{\theta}}_n)$ or the empirical variance matrix

$$\widehat{\boldsymbol{\Sigma}}_n = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{Z}_i - \bar{\mathbf{Z}}_n) (\mathbf{Z}_i - \bar{\mathbf{Z}}_n)^\top,$$

with $\mathbf{Z}_i = (t_1(\mathbf{X}_i), \dots, t_p(\mathbf{X}_i))^\top$.

Confidence intervals for θ_{Xj}

Let θ_{Xj} stand for the j -th component of the true value of the parameter $\boldsymbol{\theta}_X = (\theta_{X1}, \dots, \theta_{Xp})^\top$.

Put $\widehat{\boldsymbol{\theta}}_n = (\widehat{\theta}_{n1}, \dots, \widehat{\theta}_{np})^\top$ and $\boldsymbol{\theta}_X = (\theta_{X1}, \dots, \theta_{Xp})^\top$. By Theorem 4 we know that

$$\sqrt{n} (\widehat{\theta}_{nj} - \theta_{Xj}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, v_{jj}(\boldsymbol{\theta}_X)), \quad j \in \{1, \dots, p\},$$

where $v_{jj}(\boldsymbol{\theta}_X)$ is the j -th diagonal element of the asymptotic variance matrix

$$\mathbb{V} = \mathbb{D}_{\boldsymbol{\tau}}^{-1}(\boldsymbol{\theta}_X) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{D}_{\boldsymbol{\tau}}^{-1}(\boldsymbol{\theta}_X)]^\top. \quad (7)$$

Thus the (asymptotic two-sided) confidence interval for θ_{Xj} is given by

$$\left(\widehat{\theta}_{nj} - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\widehat{v}_{jj}}{n}}, \widehat{\theta}_{nj} + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\widehat{v}_{jj}}{n}} \right),$$

where \widehat{v}_{jj} is the j -th diagonal element of the estimated variance matrix

$$\widehat{\mathbb{V}}_n = \mathbb{D}_{\boldsymbol{\tau}}^{-1}(\widehat{\boldsymbol{\theta}}_n) \widehat{\boldsymbol{\Sigma}}_n [\mathbb{D}_{\boldsymbol{\tau}}^{-1}(\widehat{\boldsymbol{\theta}}_n)]^\top.$$

Applications of moment estimators

As maximum likelihood estimators are preferred over moment estimators, the use of moment estimators is limited. Nevertheless the moment estimators can be of interest in the following situations:

- the calculation of the maximum likelihood estimate is computationally too prohibitive due to a very complex model or a huge amount of data;
- moment estimates can be used as the starting values for the numerical algorithms that search for maximum likelihood estimates.

The choice of the functions t_1, \dots, t_p

The most common choice $t_j(x) = x^j$, where $j \in \{1, \dots, p\}$ for the univariate observations is not necessarily the most appropriate one. The idea is to choose the functions t_1, \dots, t_p so that the asymptotic variance matrix (7) is in some sense ‘minimized’. But this is usually a too difficult problem. Nevertheless one should at least check that the vector function $\tau : \Theta \rightarrow \mathbb{R}^p$ is one-to-one, otherwise the parameter θ_X might not be identifiable with the given t_1, \dots, t_p .

Now the continuity of τ guarantees the consistency of the estimator $\hat{\theta}_n$. To guarantee also the asymptotic normality one needs that the Jacobi matrix $D_\tau(\theta)$ is regular for each $\theta \in \Theta$.

To be more specific, consider the one-dimensional parameter θ and for a given function t introduce

$$\tau(\theta) = \mathbf{E}_\theta t(\mathbf{X}_1).$$

Then we need that $\tau : \Theta \rightarrow \mathbb{R}$ is a one-to-one function. Otherwise it might happen that with probability going to one the estimating function

$$\tau(\hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n t(\mathbf{X}_i).$$

has more roots (whose values are in the parameter space Θ) and we do not know which of the root is the appropriate (consistent) one.

Example 10. Let X_1, \dots, X_n be independent identically distributed random variables from the discrete distribution given as

$$\mathbf{P}(X_1 = -1) = p, \quad \mathbf{P}(X_1 = 0) = 1 - p - p^2, \quad \mathbf{P}(X_1 = 2) = p^2,$$

where $p \in \Theta = (0, \frac{-1+\sqrt{5}}{2})$.

Now the standard choice $t(x) = x$ yields that $\tau(p) = \mathbf{E}_p X_1 = 2p^2 - p$. Note that the estimating equation given by

$$2\hat{p}_n^2 - \hat{p}_n = \bar{X}_n$$

has two roots

$$\widehat{p}_n^{(1,2)} = \frac{1}{4} \pm \sqrt{\frac{\bar{X}_n}{2} + \frac{1}{16}}.$$

Show that if the true value of the parameter $p_X \in (0, \frac{1}{2})$, then

$$\widehat{p}_n^{(1)} \xrightarrow[n \rightarrow \infty]{P} \frac{1}{4} - |p_X - \frac{1}{4}|, \quad \widehat{p}_n^{(2)} \xrightarrow[n \rightarrow \infty]{P} \frac{1}{4} + |p_X - \frac{1}{4}|.$$

Thus except for the $p_X = \frac{1}{4}$ the roots $\widehat{p}_n^{(1)}$ and $\widehat{p}_n^{(2)}$ converge in distribution to different limits and only one of these limits is the true value of the parameter p_X . Note also $p_X = \frac{1}{4}$, then both the roots are consistent, but as $\tau'(\frac{1}{4}) = 0$ neither of the roots is asymptotically normal.

Show that taking $t(x) = x^2$ or simply $t(x) = \mathbb{1}\{x = -1\}$ does not introduce such problematic issues.

1.4 Confidence intervals and asymptotic variance-stabilising transformation

In this section* we are interested in constructing a confidence interval for (one-dimensional) parameter θ_X . Suppose we have an estimator $\widehat{\theta}_n$ of parameter θ_X such that

$$\sqrt{n} (\widehat{\theta}_n - \theta_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, \sigma^2(\theta_X)), \quad (8)$$

where $\sigma^2(\cdot)$ is a function continuous in the true value of the parameter (θ_X).

Standard asymptotic confidence interval of ‘Wald’ type

This interval is based on the fact that

$$\frac{\sqrt{n} (\widehat{\theta}_n - \theta_X)}{\sigma(\widehat{\theta}_n)} \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, 1)$$

and thus

$$\left(\widehat{\theta}_n - \frac{u_{1-\alpha/2} \sigma(\widehat{\theta}_n)}{\sqrt{n}}, \widehat{\theta}_n + \frac{u_{1-\alpha/2} \sigma(\widehat{\theta}_n)}{\sqrt{n}} \right) \quad (9)$$

is a confidence interval for parameter θ_X with the asymptotic coverage $1 - \alpha$.

The advantage of the confidence interval (9) is that it is easy to calculate. On the other hand the simulations show that for small sample size and/or if $|\sigma'(\theta)|$ is large then the actual coverage of this confidence interval can be much smaller than $1 - \alpha$.

* Not presented at the lecture. It is assumed that this is known from the bachelor degree.

Implicit (asymptotic) confidence interval of ‘Wilson’ type

This interval is based directly on (8) and it is given implicitly by

$$\left\{ \theta : \left| \frac{\sqrt{n}(\hat{\theta}_n - \theta)}{\sigma(\theta)} \right| < u_{1-\alpha/2} \right\}. \quad (10)$$

Note that (10) can be viewed as the set of θ for which we do not reject the null hypothesis

$$H_0 : \theta_X = \theta \quad \text{against the alternative} \quad H_1 : \theta_X \neq \theta$$

with the critical region

$$\left| \frac{\sqrt{n}(\hat{\theta}_n - \theta)}{\sigma(\theta)} \right| \geq u_{1-\alpha/2}.$$

In fact the set given by (10) does not have to be necessarily an interval. But usually the function $\theta \mapsto \frac{\sqrt{n}(\hat{\theta}_n - \theta)}{\sigma(\theta)}$ is not increasing which guarantees that the set (10) is indeed an interval.

It was observed that usually the actual coverage of this implicit confidence interval is closer to $1 - \alpha$ than for the standard asymptotic confidence interval (9). In particular if one is interested in two-sided intervals then the implicit confidence interval (10) works surprisingly well even for very small samples. Its disadvantage is that in general one does not have an explicit formula for this interval and often it has to be found with the help of methods of numerical mathematics.

Confidence interval based on the transformation stabilizing the asymptotic variance

Put $g(\theta) = \int \frac{1}{\sigma(\theta)} d\theta$. Then with the help of (8) and Δ -theorem it holds

$$\sqrt{n} (g(\hat{\theta}_n) - g(\theta_X)) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, 1).$$

Thus the set $\left(g(\hat{\theta}_n) - \frac{u_{1-\alpha/2}}{\sqrt{n}}, g(\hat{\theta}_n) + \frac{u_{1-\alpha/2}}{\sqrt{n}} \right)$ is a confidence set for $g(\theta_X)$. Now as g is an increasing function (note that $g'(\theta) > 0$) one can conclude that

$$\left(g^{-1} \left(g(\hat{\theta}_n) - \frac{u_{1-\alpha/2}}{\sqrt{n}} \right), g^{-1} \left(g(\hat{\theta}_n) + \frac{u_{1-\alpha/2}}{\sqrt{n}} \right) \right) \quad (11)$$

is a confidence interval for the parameter θ_X with the asymptotic coverage $1 - \alpha$.

The actual coverage of this confidence interval is also usually closer to $1 - \alpha$ than for the standard confidence interval (9). On the other hand when one is interested in two-sided confidence interval then the implicit confidence interval (10) usually works better. But the advantage of (11) is that one usually has an explicit formula for the confidence interval (provided that g and g^{-1} can be explicitly calculated). The confidence interval (11) is also usually a better choice than the the implicit confidence interval when one is interested in one-sided confidence intervals.

Example 11. A random sample from Poisson distribution. Find the transformation that stabilises the asymptotic variance of \bar{X}_n and based on this transformation derive the asymptotic confidence intervals for λ .

Example 12. Fisher's Z-transformation and various confidence intervals for the correlation coefficient.

Example 13. Consider a random sample from Bernoulli distribution. Find the asymptotic variance-stabilizing transformation for \bar{X}_n and construct the confidence interval based on this transformation.

Literature: [van der Vaart \[2000\]](#) – Chapters 2.1, 2.2, 3.1, 3.2 and 4.1. In particular Theorems 2.3, 2.4, 2.8 and 3.1.

The expected
end of class 3
(27. 2. 2024)

2 Maximum likelihood methods

Suppose we have a random sample of random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ being distributed as the generic vector $\mathbf{X} = (X_1, \dots, X_k)^\top$ that has a density $f(\mathbf{x}; \boldsymbol{\theta})$ with respect to a σ -finite measure μ and that the density is known up to an unknown p -dimensional parameter $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^\top \in \Theta$. Let $\boldsymbol{\theta}_X = (\theta_{X_1}, \dots, \theta_{X_p})^\top$ be the true value of the parameter.

Define the *likelihood function* as

$$L_n(\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{X}_i; \boldsymbol{\theta})$$

and the *log-likelihood function* as

$$\ell_n(\boldsymbol{\theta}) = \log L_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(\mathbf{X}_i; \boldsymbol{\theta}).$$

The *maximum likelihood estimator* of parameter $\boldsymbol{\theta}_X$ is defined as

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} L_n(\boldsymbol{\theta}) \quad \text{or alternatively as} \quad \hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}). \quad (12)$$

The (exact) distribution of $\hat{\boldsymbol{\theta}}_n$ is usually too difficult or even impossible to calculate. Thus to make the inference about $\boldsymbol{\theta}_X$ we need to derive the asymptotic distribution of $\hat{\boldsymbol{\theta}}_n$.

2.1 Asymptotic normality of maximum likelihood estimator

Regularity assumptions

Let $I(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}} \left[\frac{\partial \log f(\mathbf{X}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \frac{\partial \log f(\mathbf{X}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right]$ be the Fisher information matrix.

[R0] For any $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta$ it holds that $f(\mathbf{x}; \boldsymbol{\theta}_1) = f(\mathbf{x}; \boldsymbol{\theta}_2)$ μ -almost everywhere if and only if $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$. (*Identifiability*)

[R1] The number of parameters p in the model is constant.

[R2] The support set $S = \{\mathbf{x} \in \mathbb{R}^k : f(\mathbf{x}; \boldsymbol{\theta}) > 0\}$ does not depend on the value of the parameter $\boldsymbol{\theta}$.

[R3] (The true value of the parameter) $\boldsymbol{\theta}_X$ is an interior point of the parameter space Θ .

[R4] The density $f(\mathbf{x}; \boldsymbol{\theta})$ is three-times differentiable with respect to $\boldsymbol{\theta}$ on an open neighbourhood U of $\boldsymbol{\theta}_X$ (for μ -almost all \mathbf{x}). Further there exists a function $M(\mathbf{x})$ such that for each $j, k, l \in \{1, \dots, p\}$

$$\sup_{\boldsymbol{\theta} \in U} \left| \frac{\partial^3 \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j \partial \theta_k \partial \theta_l} \right| \leq M(\mathbf{x}),$$

for μ -almost all \mathbf{x} and

$$\mathbb{E}_{\theta_X} M(\mathbf{X}_1) < \infty.$$

[R5] The Fisher information matrix $I(\theta_X)$ is finite and positive definite.

[R6] The order of differentiation and integration can be interchanged in expressions such as

$$\frac{\partial}{\partial \theta_j} \int h(\mathbf{x}; \boldsymbol{\theta}) d\mu(\mathbf{x}) = \int \frac{\partial}{\partial \theta_j} h(\mathbf{x}; \boldsymbol{\theta}) d\mu(\mathbf{x}),$$

where $h(\mathbf{x}; \boldsymbol{\theta})$ is either $f(\mathbf{x}; \boldsymbol{\theta})$ or $\partial f(\mathbf{x}; \boldsymbol{\theta})/\partial \theta_k$ and $j, k \in \{1, \dots, p\}$.

Note that thanks to assumption **[R6]** one can calculate the Fisher information matrix as

$$I(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2 \log f(\mathbf{X}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right],$$

see for instance Lemma 5.3 of [Lehmann and Casella \[1998\]](#) or Theorem 7.27 of [Anděl \[2007\]](#).

Example 14. Let X_1, \dots, X_n be a random sample from the normal distribution $\mathbf{N}(\mu_1 + \mu_2, 1)$. Then the identifiability assumption **[R0]** is not satisfied for the vector parameter $\boldsymbol{\theta} = (\mu_1, \mu_2)^\top$.

Example 15. Let X_1, \dots, X_n be a random sample from the uniform distribution $\mathbf{U}(0, \theta)$. Note that assumption **[R2]** is not satisfied.

Show that the maximum likelihood estimator of θ is $\hat{\theta}_n = \max_{1 \leq i \leq n} \{X_i\}$. Derive the asymptotic distribution of $n(\hat{\theta}_n - \theta)$.

Remark 5. Note that in particular assumption **[R4]** is rather strict. There are ways how to derive the asymptotic normality of the maximum likelihood estimator under less strict assumptions but that would require concepts that are out of the scope of this course.

The *score function* of the i -th observation \mathbf{X}_i for the parameter $\boldsymbol{\theta}$ is defined as

$$\mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}) = \frac{\partial \log f(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

The random vector

$$\mathbf{U}_n(\boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial \log f(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

is called *the score statistic*.

We search for the maximum likelihood estimator $\hat{\boldsymbol{\theta}}_n$ as a solution of the system of the likelihood equations

$$\mathbf{U}_n(\hat{\boldsymbol{\theta}}_n) \stackrel{!}{=} \mathbf{0}_p. \quad (13)$$

Further define the observed (empirical) information matrix as

$$I_n(\boldsymbol{\theta}) = -\frac{1}{n} \frac{\partial \mathbf{U}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} = \frac{1}{n} \sum_{i=1}^n I(\mathbf{X}_i; \boldsymbol{\theta}),$$

where

$$I(\mathbf{X}_i; \boldsymbol{\theta}) = -\frac{\partial \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} = -\frac{\partial^2 \log f(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

is the contribution of the i -th observation to the information matrix.

In what follows it will be useful to prove that $I_n(\widehat{\boldsymbol{\theta}}_n) \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_X) = \mathbb{E} I(\mathbf{X}_1; \boldsymbol{\theta})$ (provided that $\widehat{\boldsymbol{\theta}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$). The following technical lemma is a generalization of this result that will be convenient in the proofs of the several theorems that will follow.

Lemma 1. *Suppose that assumptions **[R0]**-**[R6]** hold. Let ε_n be a sequence of positive numbers going to zero. Then*

$$\max_{j,k \in \{1, \dots, p\}} \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (I_n(\boldsymbol{\theta}) - I(\boldsymbol{\theta}_X))_{jk} \right| = o_P(1),$$

where

$$U_{\varepsilon_n} = \{ \boldsymbol{\theta} \in \Theta : \|\boldsymbol{\theta} - \boldsymbol{\theta}_X\| \leq \varepsilon_n \}$$

and $(I_n(\boldsymbol{\theta}) - I(\boldsymbol{\theta}_X))_{jk}$ stands for the (j, k) -element of the difference of the matrices $I_n(\boldsymbol{\theta}) - I(\boldsymbol{\theta}_X)$.

Proof. Using assumption **[R4]** and the law of large numbers one can bound

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (I_n(\boldsymbol{\theta}) - I(\boldsymbol{\theta}_X))_{jk} \right| &\leq \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (I_n(\boldsymbol{\theta}) - I_n(\boldsymbol{\theta}_X))_{jk} \right| + \left| (I_n(\boldsymbol{\theta}_X) - I(\boldsymbol{\theta}_X))_{jk} \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^p M(\mathbf{X}_i) \varepsilon_n + o_P(1) = O_P(1) o(1) + o_P(1) = o_P(1), \end{aligned}$$

which implies the statement of the lemma. □

Corollary 1. *Let the assumptions of Lemma 1 be satisfied. Further let $\widehat{\mathbf{t}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$. Then for each $j, k \in \{1, \dots, p\}$*

$$\left| (I_n(\widehat{\mathbf{t}}_n) - I(\boldsymbol{\theta}_X))_{jk} \right| = o_P(1).$$

Proof. Note that $\widehat{\mathbf{t}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$ implies that there exists a sequence of positive constants $\{\varepsilon_n\}$ going to zero such that

$$\mathbb{P}(\widehat{\mathbf{t}}_n \in U_{\varepsilon_n}) \xrightarrow[n \rightarrow \infty]{} 1.$$

The corollary now follows from Lemma 1 and from the fact that one can bound

$$\begin{aligned} \left| (I_n(\widehat{\mathbf{t}}_n) - I(\boldsymbol{\theta}_X))_{jk} \right| &= \left| (I_n(\widehat{\mathbf{t}}_n) - I(\boldsymbol{\theta}_X))_{jk} \right| \mathbb{1}\{\widehat{\mathbf{t}}_n \in U_{\varepsilon_n}\} + \left| (I_n(\widehat{\mathbf{t}}_n) - I(\boldsymbol{\theta}_X))_{jk} \right| \mathbb{1}\{\widehat{\mathbf{t}}_n \notin U_{\varepsilon_n}\} \\ &\leq \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (I_n(\boldsymbol{\theta}) - I(\boldsymbol{\theta}_X))_{jk} \right| + \left| (I_n(\widehat{\mathbf{t}}_n) - I(\boldsymbol{\theta}_X))_{jk} \right| \mathbb{1}\{\widehat{\mathbf{t}}_n \notin U_{\varepsilon_n}\}. \end{aligned}$$

□

Theorem 5. *Suppose that assumptions [R0]-[R6] hold.*

(i) *Then with probability tending to one as $n \rightarrow \infty$ there exists a consistent solution $\widehat{\boldsymbol{\theta}}_n$ of the likelihood equations (13).**

(ii) *Any consistent solution $\widehat{\boldsymbol{\theta}}_n$ of the likelihood equations (13) satisfies,*

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = [I(\boldsymbol{\theta}_X)]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) + o_P(1), \quad (14)$$

which further implies that

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, I^{-1}(\boldsymbol{\theta}_X)). \quad (15)$$

Proof of (i). First, we need to prove the existence of the consistent root $\widehat{\boldsymbol{\theta}}_n$ of the likelihood equations. This can be deduced from a more general Theorem 9. An alternative approach can be found in the proof of Theorem 5.1 of Lehmann and Casella [1998, Chapter 6].

Proof of (ii). Suppose that $\widehat{\boldsymbol{\theta}}_n$ is a consistent solution of the likelihood equations. Then by the mean value theorem (applied to each component of $\mathbf{U}_n(\boldsymbol{\theta})$) one gets that

$$\mathbf{0}_p = \mathbf{U}_n(\widehat{\boldsymbol{\theta}}_n) = \mathbf{U}_n(\boldsymbol{\theta}_X) - n I_n^*(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X),$$

where I_n^* is a matrix with the elements

$$i_{n,jk}^* = \frac{1}{n} \sum_{i=1}^n \frac{-\partial^2 \log f(\mathbf{X}_i; \boldsymbol{\theta})}{\partial \theta_j \partial \theta_k} \Big|_{\boldsymbol{\theta} = \widehat{\mathbf{t}}_n^{(j)}}, \quad j, k \in \{1, \dots, p\},$$

with $\widehat{\mathbf{t}}_n^{(j)}$ being between $\widehat{\boldsymbol{\theta}}_n$ and $\boldsymbol{\theta}_X$. Thus the consistency of $\widehat{\boldsymbol{\theta}}_n$ implies that $\widehat{\mathbf{t}}_n^{(j)} \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$ and one can use Corollary 1 to show that

$$I_n^* \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_X). \quad (16)$$

Thus with probability going to one there exists $[I_n^*]^{-1}$ and one can write

$$n(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = [I_n^*]^{-1} \mathbf{U}_n(\boldsymbol{\theta}_X).$$

* Thus defining the estimator as an appropriately chosen root of the likelihood equations (provided that the likelihood equations has at least one root) and zero otherwise yields a consistent estimator of $\boldsymbol{\theta}_X$.

Now the central limit theorem for independent identically distributed random vectors implies that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, I(\boldsymbol{\theta}_X)). \quad (17)$$

Note that (17) yields that $\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) = O_P(1)$. Thus using (16) and CMT (Theorem 1) implies that

$$\begin{aligned} \sqrt{n} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) &= [I_n^*]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) \\ &= [I^{-1}(\boldsymbol{\theta}_X) + o_P(1)] \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) \\ &= I^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \boldsymbol{\theta}_X) + o_P(1). \end{aligned}$$

Now (15) follows by CS (Theorem 2) and (17). □

Remark 6. While the proof of consistency is for $p = 1$ relatively simple [see e.g. Theorem 22 of Nagy], for $p > 1$ it is much more involved. The reason is that while the border of the neighbourhood in \mathbb{R} is a two-point set, in \mathbb{R}^p ($p > 1$) it is an uncountable set.

Remark 7. Note that strictly speaking Theorem 5 does not guarantee the asymptotic normality of the maximum likelihood estimator but of an appropriately chosen root of the likelihood equations (13). As illustrated in Example 19 it may happen that the maximum likelihood estimator defined by (12) is not a consistent estimator of $\boldsymbol{\theta}_X$ even if all the regularity assumptions [R0]-[R6] are satisfied. It may also happen that the maximum likelihood estimator does not exist (see the example on page 21). That is why some authors define the maximum likelihood estimator in regular families as an appropriately chosen root of the likelihood equations.

Fortunately for many models commonly used in applications the log-likelihood function $\ell_n(\boldsymbol{\theta})$ is (almost surely) convex. Then the maximum likelihood estimator is the only solution to the likelihood equations and Theorem 5 guarantees that this estimator is asymptotically normal. If $\ell_n(\boldsymbol{\theta})$ is not convex, there might be more roots to the likelihood equations and the choice of an appropriate (consistent) root of the estimating equations is more delicate both from the theoretical as well as the numerical point of view. Other available consistent estimators (e.g. moment estimators) can be very useful as for instance the starting points of the numerical algorithms that search for the root of the likelihood equations.

Example 16. Let X_1, \dots, X_n be a random sample from Bernoulli distribution* $\text{Be}(p)$. Note that if either $\sum_{i=1}^n X_i = 0$ or $\sum_{i=1}^n X_i = n$ then there is no root of the likelihood equation. Nevertheless the probability of both events converges to zero as $n \rightarrow \infty$ whenever $p_X \in (0, 1)$.

* Alternativního rozdělení

Example 17. Let X_1, \dots, X_n be a random sample from the Pareto distribution with the density

$$f(x) = \frac{\beta \alpha^\beta}{x^{\beta+1}} \mathbb{1}\{x \geq \alpha\}, \quad \beta > 0, \alpha > 0,$$

where both parameters are unknown.

- (i) Find the maximum likelihood estimator of $\hat{\boldsymbol{\theta}}_n = (\hat{\alpha}_n, \hat{\beta}_n)^\top$ of the parameter $\boldsymbol{\theta} = (\alpha, \beta)^\top$.
- (ii) Derive the asymptotic distribution of $n(\hat{\alpha}_n - \alpha)$.
- (iii) Derive the asymptotic distribution of $\sqrt{n}(\hat{\beta}_n - \beta)$.

Example 18. Let X_1, \dots, X_n be a random sample from $\mathbf{N}(\mu, 1)$ where the parameter space for the parameter μ is restricted to $[0, \infty)$. Find the maximum likelihood estimator of μ and derive its asymptotic distribution. Do not forget to consider the special case $\mu = 0$.

Example 19. Let X_1, \dots, X_n be a random sample from the mixture of distributions $\mathbf{N}(0, 1)$ and $\mathbf{N}(\theta, \exp\{-2/\theta^2\})$ with equal weights and the parameter space given by $\Theta = (0, \infty)$. Define the estimator of the parameter θ as $\hat{\theta}_n^{(ML)} = \arg \max_{\theta \in \Theta} \ell_n(\theta)$. Then it can be shown that $\hat{\theta}_n^{(ML)} \xrightarrow[n \rightarrow \infty]{P} 0$, thus $\hat{\theta}_n^{(ML)}$ is not consistent estimator.

Nevertheless note that the assumptions **[R0]**-**[R6]** are met. Thus by Theorem 5 there exists a different root ($\hat{\theta}_n$) of the likelihood equation such that this estimator satisfies (14) and (15).

Example 20. Let X_1, \dots, X_n be a random sample from the mixture of distributions $\mathbf{N}(0, 1)$ and $\mathbf{N}(\mu, \sigma^2)$ with equal weights and the parameter space for the parameter $\boldsymbol{\theta} = (\mu, \sigma)^\top$ is given by $\Theta = \mathbb{R} \times (0, \infty)$. Show that

$$\sup_{(\mu, \sigma^2)^\top \in \Theta} \ell_n(\mu, \sigma^2) = \infty$$

and that the maximum likelihood estimator does not exist. But similarly as in Example 19 Theorem 5 still holds.

2.2 Asymptotic efficiency of maximum likelihood estimators

Recall the Rao-Cramér inequality. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from the regular family of densities $\mathcal{F} = \{f(\mathbf{x}; \boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta\}$, and \mathbf{T}_n be an *unbiased* estimator of $\boldsymbol{\theta}_X$ (based on $\mathbf{X}_1, \dots, \mathbf{X}_n$). Then

$$\text{var}(\mathbf{T}_n) - \frac{1}{n} I^{-1}(\boldsymbol{\theta}_X) \geq 0.$$

By Theorem 5 we have that (under appropriate regularity assumptions)

$$\text{avar}(\hat{\boldsymbol{\theta}}_n) = \frac{1}{n} I^{-1}(\boldsymbol{\theta}_X).$$

Thus the asymptotic variance of $\hat{\boldsymbol{\theta}}_n$ attains the lower bound in Rao-Cramér inequality.

Remark 8. Note that strictly speaking comparing with the Rao-Cramér bound is not fair. Generally, the maximum likelihood estimator $\widehat{\boldsymbol{\theta}}_n$ is not unbiased. Further, Rao-Cramér inequality speaks about the bound on the variance, but we compare the asymptotic variance of $\widehat{\boldsymbol{\theta}}_n$ with this bound. Nevertheless it can be shown that in regular models there exists a lower bound for the asymptotic variances of the estimators that are asymptotically normal with zero mean and in some (natural) sense regular (see Example 21 below). And this bound is indeed given by $\frac{1}{n} I^{-1}(\boldsymbol{\theta}_X)$. See also Serfling [1980, Chapter 4.1.3] and the references therein.

Example 21. Let X_1, \dots, X_n be a random sample from $\mathbf{N}(\theta, 1)$, where $\theta \in \mathbb{R}$. Define the estimator of θ as

$$\widehat{\theta}_n^{(S)} = \begin{cases} 0, & \text{if } |\bar{X}_n| \leq n^{-1/4}, \\ \bar{X}_n, & \text{if } |\bar{X}_n| > n^{-1/4}. \end{cases}$$

This estimator is called also *Hodges* or *shrinkage* estimator. Show that if $\theta_X \neq 0$ then $\sqrt{n}(\widehat{\theta}_n^{(S)} - \theta_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, 1)$ and if $\theta_X = 0$ then even $n^r(\widehat{\theta}_n^{(S)} - \theta_X) \xrightarrow[n \rightarrow \infty]{P} 0$ for each $r \in \mathbb{N}$. Thus from the point-wise asymptotic point of view, the estimator $\widehat{\theta}_n^{(S)}$ is better than the standard maximum likelihood estimator that is given by the sample mean \bar{X}_n .

But on the other hand consider the following sequence of the true values of the parameter $\theta_X^{(n)} = n^{-1/4}$. Then show that for an arbitrarily large value of K

$$\liminf_{n \rightarrow \infty} \mathbf{P} \left(\sqrt{n} (\widehat{\theta}_n^{(S)} - \theta_X^{(n)}) \geq K \right) \geq \frac{1}{2}.$$

Thus the sequence $\sqrt{n}(\widehat{\theta}_n^{(S)} - \theta_X^{(n)})$ is not tight and so it does not converge in distribution. Such a non-uniform behaviour of the estimator $\widehat{\theta}_n^{(S)}$ is usually considered as undesirable. Thus the aim of the regularity assumptions on the estimators is to avoid such estimators that from the point-wise view can be considered as superior (*superefficient*) to the maximum likelihood estimators.*

The end of
class 5
(5. 3. 2024)

2.3 Estimation of the asymptotic variance matrix

To do the inference about the parameter $\boldsymbol{\theta}_X$ we need to have a consistent estimator of $I(\boldsymbol{\theta}_X)$. Usually, we use one of the following estimators

$$I(\widehat{\boldsymbol{\theta}}_n) \quad \text{or} \quad I_n(\widehat{\boldsymbol{\theta}}_n) \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) \mathbf{U}^\top(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n).$$

The consistency of $I(\widehat{\boldsymbol{\theta}}_n)$ follows by CMT (Theorem 1), provided (the matrix function) $I(\boldsymbol{\theta})$ is continuous in $\boldsymbol{\theta}_X$, which follows by assumption [R4].

* Note that the issue of superefficiency is behind the claimed ‘oracle’-properties of some regularized estimators (e.g. adaptive LASSO), see Leeb and Pötscher [2008] and the references therein.

The consistency of $I_n(\widehat{\boldsymbol{\theta}}_n) \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_X)$ follows from Corollary 1 and Theorem 5.

On the other hand the consistency of $\frac{1}{n} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) \mathbf{U}^\top(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n)$ does not automatically follow from assumptions [R0]-[R6]. It can be proved analogously as Corollary 1 provided the following assumption holds.

[R7] There exists an open neighbourhood U of $\boldsymbol{\theta}_X$ such that for each j, k in $\{1, \dots, p\}$ there exists a function $M_{jkl}(\mathbf{x})$ such that

$$\sup_{\boldsymbol{\theta} \in U} \left| \frac{\partial^2 \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j \partial \theta_k} \right| \leq M_{jk}(\mathbf{x})$$

for μ -almost all \mathbf{x} and

$$\mathbb{E}_{\boldsymbol{\theta}_X} M_{jk}^2(\mathbf{X}_i) < \infty.$$

Literature: Anděl [2007] Chapter 7.6.5, Lehmann and Casella [1998] Chapter 6.5, Kulich [2014].

2.4 Asymptotic tests (without nuisance parameters)

Suppose we are interested in testing the null hypothesis

$$H_0 : \boldsymbol{\theta}_X = \boldsymbol{\theta}_0 \text{ against the alternative } H_1 : \boldsymbol{\theta}_X \neq \boldsymbol{\theta}_0.$$

Let \widehat{I}_n be an estimate of the Fisher information matrix $I(\boldsymbol{\theta}_X)$ or $I(\boldsymbol{\theta}_0)$. Basically there are three tests that can be considered.

Likelihood ratio test is based on the test statistic

$$LR_n = 2 (\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\boldsymbol{\theta}_0)).$$

Wald test is based on the test statistic

$$W_n = n (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^\top \widehat{I}_n (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0).$$

Rao score test is based on the test statistic

$$R_n = \frac{1}{n} \mathbf{U}_n^\top(\boldsymbol{\theta}_0) \widehat{I}_n^{-1} \mathbf{U}_n(\boldsymbol{\theta}_0). \quad (18)$$

Note that the advantage of the likelihood ratio test (LR_n) is that one does not need to estimate the Fisher information matrix. On the other hand the advantage of Rao score test (R_n) is that you do not need to calculate the maximal likelihood estimator $\widehat{\boldsymbol{\theta}}_n$. That is why in Rao score statistic (R_n) one uses usually either $I(\boldsymbol{\theta}_0)$ or $I_n(\boldsymbol{\theta}_0)$ as \widehat{I}_n . On the other hand usually (for historical reasons) $I(\widehat{\boldsymbol{\theta}}_n)$ or $I_n(\widehat{\boldsymbol{\theta}}_n)$ is used for Wald statistic (W_n).

The next theorem says that all the test statistics have the same asymptotic distribution under the null hypothesis.

Theorem 6. Suppose that the null hypothesis holds, assumptions **[R0]**-**[R6]** are satisfied, $\widehat{I}_n \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_0)$ and $\widehat{\boldsymbol{\theta}}_n$ is a consistent solution of the likelihood equations. Then each of the test statistics LR_n , W_n and R_n converges in distribution to χ^2 -distribution with p degrees of freedom.

Proof. R_n : Note that R_n can be rewritten as

$$R_n = \left([\widehat{I}_n]^{-\frac{1}{2}} \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_0) \right)^\top \left([\widehat{I}_n]^{-\frac{1}{2}} \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_0) \right).$$

Now by the asymptotic normality of the score statistic (17), consistency of \widehat{I}_n and CS (Theorem 2) one gets that

$$[\widehat{I}_n]^{-\frac{1}{2}} \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{I}_p),$$

where \mathbb{I}_p is an identity matrix of dimension $p \times p$. Now the statement follows by using CMT (Theorem 1) with $g(x_1, \dots, x_p) = \sum_{j=1}^p x_j^2$.

W_n : One can rewrite W_n as

$$W_n = \left([\widehat{I}_n]^{\frac{1}{2}} \sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \right)^\top \left([\widehat{I}_n]^{\frac{1}{2}} \sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \right).$$

Now the statement follows by analogous reasoning as for R_n , as by Theorem 5 and CS (Theorem 2) one gets

$$[\widehat{I}_n]^{\frac{1}{2}} \sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{I}_p).$$

LR_n : With the help of the second order Taylor expansion around $\widehat{\boldsymbol{\theta}}_n$ one gets:

$$\ell_n(\boldsymbol{\theta}_0) = \ell_n(\widehat{\boldsymbol{\theta}}_n) + \underbrace{\mathbf{U}_n^\top(\widehat{\boldsymbol{\theta}}_n)}_{=\mathbf{0}_p^\top} (\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}_n) - \frac{n}{2} (\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}_n)^\top I_n(\boldsymbol{\theta}_n^*) (\boldsymbol{\theta}_0 - \widehat{\boldsymbol{\theta}}_n),$$

where $\boldsymbol{\theta}_n^*$ lies between $\boldsymbol{\theta}_0$ and $\widehat{\boldsymbol{\theta}}_n$. Applying Corollary 1 yields $I_n(\boldsymbol{\theta}_n^*) \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_0)$. Thus analogously as above one gets

$$LR_n = 2(\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\boldsymbol{\theta}_0)) = \sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^\top I_n(\boldsymbol{\theta}_n^*) \sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow[n \rightarrow \infty]{d} \chi_p^2.$$

□

Remark 9. Note that using the asymptotic representation (14) of the maximum likelihood estimator $\widehat{\boldsymbol{\theta}}_n$ and the derivations done in the proof of Theorem 6 one can show that the difference of each of the two test statistics (LR_n , W_n and R_n) converges under the null hypothesis to zero in probability.

Nevertheless, in simulations it is observed that the actual level (the probability of type one error) of the test for the Wald test (W_n) can be substantially different from the prescribed

level α . Unfortunately, usually the test is anti-conservative, i.e. the actual level is higher than the prescribed level α . This happens in particular for small samples and/or when the curvature of the log-likelihood $\ell_n(\boldsymbol{\theta})$ is relatively high (as measured for instance by $I(\boldsymbol{\theta})$). The latter happens often if $\boldsymbol{\theta}_0$ is close to the border of the parameter space Θ . That is why some authors recommend either the score test R_n or likelihood ratio test LR_n whose actual levels are usually very close to the prescribed level α even in small samples.

Example 22. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample of K -variate random vectors from the multinomial distribution $\text{Mult}_K(1, \mathbf{p})$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{iK})^\top$ and $\mathbf{p} = (p_1, \dots, p_K)^\top$. Suppose we are interested in testing the null hypothesis

$$H_0 : \mathbf{p}_X = \mathbf{p}^0, \quad H_1 : \mathbf{p}_X \neq \mathbf{p}^0,$$

where $\mathbf{p}^0 = (p_1^0, \dots, p_K^0)^\top$ is a given value of the parameter \mathbf{p} . For $k \in \{1, \dots, K\}$ put $n_k = \sum_{i=1}^n X_{ik}$. Derive that

$$LR_n = 2 \sum_{k=1}^K n_k \log \left(\frac{n_k}{np_k^0} \right).$$

Further if one uses $I(\hat{\boldsymbol{\theta}}_n)$ in the Wald test and $I(\boldsymbol{\theta}_0)$ in the Rao score test, then

$$W_n = \sum_{k=1}^K \frac{(n_k - np_k^0)^2}{n_k}, \quad R_n = \sum_{k=1}^K \frac{(n_k - np_k^0)^2}{np_k^0}.$$

Show that each of the test statistics converges to χ^2 -distribution with $K - 1$ degrees of freedom.

Note that Rao score test (R_n) corresponds to the standard χ^2 -test of goodness-of-fit in multinomial distribution.

Hint. One has to be careful as it is not possible to take $\boldsymbol{\theta} = (p_1, \dots, p_K)^\top$, as $p_K = 1 - \sum_{k=1}^{K-1} p_k$ (which violates assumption **[R3]**, as the corresponding parameter space would not have any interior points). To avoid this problem one has to take for instance $\boldsymbol{\theta} = (p_1, \dots, p_{K-1})^\top$.

2.5 Asymptotic confidence sets

Sometimes we are interested in the confidence set for the whole vector parameter $\boldsymbol{\theta}_X$. Then we usually use the following confidence set

$$\left\{ \boldsymbol{\theta} \in \Theta : n (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta})^\top \hat{I}_n (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \leq \chi_p^2(1 - \alpha) \right\},$$

where \hat{I}_n is a consistent estimator of $I(\boldsymbol{\theta}_X)$. Usually $I_n(\hat{\boldsymbol{\theta}}_n)$ or $I(\hat{\boldsymbol{\theta}}_n)$ are used as \hat{I}_n . Then the resulting confidence set is an ellipsoid.

Confidence intervals for θ_{Xj}

In most of the applications we are interested in confidence intervals for components θ_{Xj} of the parameter $\boldsymbol{\theta}_X = (\theta_{X1}, \dots, \theta_{Xp})^\top$.

Put $\widehat{\boldsymbol{\theta}}_n = (\widehat{\theta}_{n1}, \dots, \widehat{\theta}_{np})^\top$ and $\boldsymbol{\theta}_X = (\theta_{X1}, \dots, \theta_{Xp})^\top$. By Theorem 5 we know that

$$\sqrt{n} (\widehat{\theta}_{nj} - \theta_{Xj}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, i^{jj}(\boldsymbol{\theta}_X)), \quad j \in \{1, \dots, p\},$$

where $i^{jj}(\boldsymbol{\theta}_X)$ is the j -th diagonal element of $I^{-1}(\boldsymbol{\theta}_X)$. Thus the asymptotic variance of $\widehat{\theta}_{jn}$ is given by $\text{avar}(\widehat{\theta}_{nj}) = \frac{i^{jj}(\boldsymbol{\theta}_X)}{n}$, which can be estimated by $\widehat{\text{avar}}(\widehat{\theta}_{nj}) = \frac{i_n^{jj}}{n}$, where i_n^{jj} is the j -th diagonal element of \widehat{I}_n^{-1} . Thus the two-sided (asymptotic) confidence interval for θ_{Xj} is given by

$$\left(\widehat{\theta}_{nj} - u_{1-\frac{\alpha}{2}} \sqrt{\frac{i_n^{jj}}{n}}, \widehat{\theta}_{nj} + u_{1-\frac{\alpha}{2}} \sqrt{\frac{i_n^{jj}}{n}} \right). \quad (19)$$

Remark 10. The approaches presented in this section are based on the Wald test statistic. The approaches based on the other test statistics are also possible. For instance one can construct the confidence set for $\boldsymbol{\theta}_X$ as

$$\{\boldsymbol{\theta} \in \Theta : 2 (\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\boldsymbol{\theta})) \leq \chi_p^2(1 - \alpha)\}.$$

But such a confidence set is for $p > 1$ very difficult to calculate. Nevertheless, as we will see later there exists an approach to calculate the confidence interval for θ_{Xj} with the help of the profile likelihood.

2.6 Asymptotic tests with nuisance parameters

Denote $\boldsymbol{\tau}$ the first q ($1 \leq q < p$) components of the vector $\boldsymbol{\theta}$ and $\boldsymbol{\psi}$ the remaining $p - q$ components, i.e.

$$\boldsymbol{\theta} = (\boldsymbol{\tau}^\top, \boldsymbol{\psi}^\top)^\top = (\theta_1, \dots, \theta_q, \theta_{q+1}, \dots, \theta_p)^\top.$$

We want to test the null hypothesis that

$$H_0 : \boldsymbol{\tau}_X = \boldsymbol{\tau}_0, \quad H_1 : \boldsymbol{\tau}_X \neq \boldsymbol{\tau}_0$$

and the remaining parameters $\boldsymbol{\psi}$ are considered as nuisance*. In regression problems this corresponds to situation when one wants to test that a given regressor (interaction) has an effect on the response. Then one is testing that all the parameters corresponding to this regressor (interaction) are zero.

In what follows all the vectors and matrices appearing in the notation of maximum likelihood estimation theory are decomposed into the first q (part 1) and the remaining $p - q$

* *rušivé*

components (part 2), i.e.

$$\widehat{\boldsymbol{\theta}}_n = \begin{pmatrix} \widehat{\boldsymbol{\tau}}_n \\ \widehat{\boldsymbol{\psi}}_n \end{pmatrix}, \quad \mathbf{U}_n(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{U}_{1n}(\boldsymbol{\theta}) \\ \mathbf{U}_{2n}(\boldsymbol{\theta}) \end{pmatrix},$$

and

$$I(\boldsymbol{\theta}) = \begin{pmatrix} I_{11}(\boldsymbol{\theta}) & I_{12}(\boldsymbol{\theta}) \\ I_{21}(\boldsymbol{\theta}) & I_{22}(\boldsymbol{\theta}) \end{pmatrix}, \quad I_n(\boldsymbol{\theta}) = \begin{pmatrix} I_{11n}(\boldsymbol{\theta}) & I_{12n}(\boldsymbol{\theta}) \\ I_{21n}(\boldsymbol{\theta}) & I_{22n}(\boldsymbol{\theta}) \end{pmatrix}. \quad (20)$$

Lemma 2. Let \mathbb{J} be a symmetric non-singular matrix of order $p \times p$ that can be written in the block form as

$$\mathbb{J} = \begin{pmatrix} \mathbb{J}_{11} & \mathbb{J}_{12} \\ \mathbb{J}_{21} & \mathbb{J}_{22} \end{pmatrix}.$$

Denote

$$\mathbb{J}_{11.2} = \mathbb{J}_{11} - \mathbb{J}_{12}\mathbb{J}_{22}^{-1}\mathbb{J}_{21}, \quad \mathbb{J}_{22.1} = \mathbb{J}_{22} - \mathbb{J}_{21}\mathbb{J}_{11}^{-1}\mathbb{J}_{12}.$$

Then

$$\mathbb{J}^{-1} = \begin{pmatrix} \mathbb{J}^{11} & \mathbb{J}^{12} \\ \mathbb{J}^{21} & \mathbb{J}^{22} \end{pmatrix},$$

where

$$\mathbb{J}^{11} = \mathbb{J}_{11.2}^{-1}, \quad \mathbb{J}^{22} = \mathbb{J}_{22.1}^{-1}, \quad \mathbb{J}^{12} = -\mathbb{J}_{11.2}^{-1}\mathbb{J}_{12}\mathbb{J}_{22}^{-1}, \quad \mathbb{J}^{21} = -\mathbb{J}_{22.1}^{-1}\mathbb{J}_{21}\mathbb{J}_{11}^{-1}.$$

Proof. Calculate $\mathbb{J}^{-1}\mathbb{J}$. □

Suppose that the parametric space can be written as $\Theta = \Theta_{\boldsymbol{\tau}} \times \Theta_{\boldsymbol{\psi}}$, where $\Theta_{\boldsymbol{\tau}} \subset \mathbb{R}^q$ and $\Theta_{\boldsymbol{\psi}} \subset \mathbb{R}^{p-q}$.

Denote $\widetilde{\boldsymbol{\theta}}_n$ the estimator of $\boldsymbol{\theta}$ under the null hypothesis, i.e.

$$\widetilde{\boldsymbol{\theta}}_n = \begin{pmatrix} \boldsymbol{\tau}_0 \\ \widetilde{\boldsymbol{\psi}}_n \end{pmatrix}, \quad \text{where } \widetilde{\boldsymbol{\psi}}_n \text{ solves } \mathbf{U}_{2n}(\boldsymbol{\tau}_0, \widetilde{\boldsymbol{\psi}}_n) \stackrel{!}{=} \mathbf{0}_{p-q}.$$

Let \widehat{I}_n^{11} be an estimate of the corresponding block $I^{11}(\boldsymbol{\theta}_X)$ in the inverse of Fisher information matrix $I^{-1}(\boldsymbol{\theta}_X)$. The three asymptotic tests of the null hypothesis $H_0 : \boldsymbol{\tau}_X = \boldsymbol{\tau}_0$ are as follows.

Likelihood ratio test is based on the test statistic

$$LR_n^* = 2 (\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\widetilde{\boldsymbol{\theta}}_n)). \quad (21)$$

Wald test is based on the test statistic

$$W_n^* = n (\widehat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0)^\top [\widehat{I}_n^{11}]^{-1} (\widehat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0).$$

Rao score test is based on the test statistic

$$R_n^* = \frac{1}{n} \mathbf{U}_{1n}^\top(\widetilde{\boldsymbol{\theta}}_n) \widehat{I}_n^{11} \mathbf{U}_{1n}(\widetilde{\boldsymbol{\theta}}_n). \quad (22)$$

The end of
class 6
(8. 3. 2024)

Remark 11. As $\mathbf{U}_{2n}(\tilde{\boldsymbol{\theta}}_n) = \mathbf{0}_{p-q}$, the test statistic of the Rao score test can be also written in a form

$$R_n^* = \frac{1}{n} \mathbf{U}_n^\top(\tilde{\boldsymbol{\theta}}_n) \hat{I}_n^{-1} \mathbf{U}_n(\tilde{\boldsymbol{\theta}}_n),$$

which is a straightforward analogy of the test statistic (18) of the Rao score test in case of no nuisance parameters.

Similarly as in the previous section the advantage of the likelihood ratio test (LR_n^*) is that one does not need to estimate $I^{-1}(\boldsymbol{\theta}_X)$. On the other hand the advantage of Rao score test (R_n^*) is that it is sufficient to calculate the maximal likelihood estimator only under the null hypothesis.

The next theorem is an analogy to Theorem 6. It says that all the test statistics have the same asymptotic distribution under the null hypothesis.

Theorem 7. *Suppose that the null hypothesis holds, assumptions [R0]-[R6] are satisfied and $\hat{I}_n^{11} \xrightarrow[n \rightarrow \infty]{P} I^{11}(\boldsymbol{\theta}_X)$. Further assume that both $\hat{\boldsymbol{\theta}}_n$ and $\tilde{\boldsymbol{\theta}}_n$ are consistent estimators of $\boldsymbol{\theta}_X$. Then each of the test statistics LR_n^* , W_n^* and R_n^* converges in distribution to χ^2 -distribution with q degrees of freedom.*

Proof. First note if the null hypothesis holds then $\boldsymbol{\theta}_X = (\boldsymbol{\tau}_0^\top, \boldsymbol{\psi}_X^\top)^\top$, where $\boldsymbol{\psi}_X$ stands for the true value of $\boldsymbol{\psi}$.

W_n^* : Note that by Theorem 5 $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, I^{-1}(\boldsymbol{\theta}_X))$, which yields

$$\sqrt{n}(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_q(\mathbf{0}_q, I^{11}(\boldsymbol{\theta}_X)).$$

Thus analogously as in the proof of Theorem 6 one can show that

$$\left[\hat{I}_n^{11} \right]^{-\frac{1}{2}} \sqrt{n}(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_q(\mathbf{0}_q, \mathbb{I}_q),$$

which further with the CMT (Theorem 1) implies

$$W_n^* = \left\{ \left[\hat{I}_n^{11} \right]^{-\frac{1}{2}} \sqrt{n}(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) \right\}^\top \left\{ \left[\hat{I}_n^{11} \right]^{-\frac{1}{2}} \sqrt{n}(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) \right\} \xrightarrow[n \rightarrow \infty]{d} \chi_q^2.$$

R_n^* : By the mean value theorem (applied to each component of $\mathbf{U}_{1n}(\boldsymbol{\theta})$) one gets

$$\frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\tilde{\boldsymbol{\theta}}_n) = \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\boldsymbol{\theta}_X) - I_{12n}^* \sqrt{n}(\tilde{\boldsymbol{\psi}}_n - \boldsymbol{\psi}_X), \quad (23)$$

where I_{12n}^* is the (1, 2)-block of the observed Fisher matrix whose j -th row ($j \in \{1, \dots, q\}$) is evaluated at some $\boldsymbol{\theta}_n^{j*}$ that is between $\tilde{\boldsymbol{\theta}}_n$ and $\boldsymbol{\theta}_X$. As $\boldsymbol{\theta}_n^{j*} \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$, Corollary 1 implies that

$$I_{12n}^* \xrightarrow[n \rightarrow \infty]{P} I_{12}(\boldsymbol{\theta}_X). \quad (24)$$

Further note that $\tilde{\boldsymbol{\psi}}_n$ is a maximum likelihood estimator in the model

$$\mathcal{F}_0 = \{f(\mathbf{x}; \boldsymbol{\tau}_0, \boldsymbol{\psi}); \boldsymbol{\psi} \in \Theta_\psi\}.$$

As the null hypothesis holds, using Theorem 5 one gets

$$\sqrt{n}(\tilde{\boldsymbol{\psi}}_n - \boldsymbol{\psi}_X) = I_{22}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) + o_P(1). \quad (25)$$

Combining (23), (24) and (25) yields

$$\frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\tilde{\boldsymbol{\theta}}_n) = \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\boldsymbol{\theta}_X) - I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) + o_P(1). \quad (26)$$

Now using (26) and the central limit theorem (for i.i.d. vectors), which implies that (written in a block form)

$$\frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X) = \begin{pmatrix} \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\boldsymbol{\theta}_X) \\ \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) \end{pmatrix} \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p \left(\mathbf{0}_p, \begin{pmatrix} I_{11}(\boldsymbol{\theta}_X) & I_{12}(\boldsymbol{\theta}_X) \\ I_{21}(\boldsymbol{\theta}_X) & I_{22}(\boldsymbol{\theta}_X) \end{pmatrix} \right),$$

one gets

$$\begin{aligned} \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\tilde{\boldsymbol{\theta}}_n) &= \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\boldsymbol{\theta}_X) - I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) + o_P(1) \\ &= (\mathbb{1}_q, -I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X)) \begin{pmatrix} \frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\boldsymbol{\theta}_X) \\ \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) \end{pmatrix} + o_P(1) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_q(0, K(\boldsymbol{\theta}_X)), \end{aligned}$$

where

$$\begin{aligned} K(\boldsymbol{\theta}_X) &= (\mathbb{1}_q, -I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X)) \begin{pmatrix} I_{11}(\boldsymbol{\theta}_X) & I_{12}(\boldsymbol{\theta}_X) \\ I_{21}(\boldsymbol{\theta}_X) & I_{22}(\boldsymbol{\theta}_X) \end{pmatrix} \begin{pmatrix} \mathbb{1}_q \\ -I_{22}^{-1}(\boldsymbol{\theta}_X) I_{21}(\boldsymbol{\theta}_X) \end{pmatrix} \\ &= I_{11}(\boldsymbol{\theta}_X) - 2I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) I_{21}(\boldsymbol{\theta}_X) + I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) I_{22}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) I_{21}(\boldsymbol{\theta}_X) \\ &= I_{11}(\boldsymbol{\theta}_X) - I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) I_{21}(\boldsymbol{\theta}_X) = I_{11:2}(\boldsymbol{\theta}_X) \stackrel{\text{Lemma 2}}{=} [I^{11}(\boldsymbol{\theta}_X)]^{-1}. \end{aligned}$$

Thus $\frac{1}{\sqrt{n}} \mathbf{U}_{1n}(\tilde{\boldsymbol{\theta}}_n) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_q(\mathbf{0}_q, [I^{11}(\boldsymbol{\theta}_X)]^{-1})$, which further with the help of CS (Theorem 2) and CMT (Theorem 1) implies the statement of the theorem for R_n^* .

LR $_n^*$: By the second-order Taylor expansion around the point $\hat{\boldsymbol{\theta}}_n$ one gets

$$\ell_n(\tilde{\boldsymbol{\theta}}_n) = \ell_n(\hat{\boldsymbol{\theta}}_n) + \underbrace{\mathbf{U}_n^\top(\hat{\boldsymbol{\theta}}_n)}_{=\mathbf{0}_p^\top} (\tilde{\boldsymbol{\theta}}_n - \hat{\boldsymbol{\theta}}_n) - \frac{n}{2} (\tilde{\boldsymbol{\theta}}_n - \hat{\boldsymbol{\theta}}_n)^\top I_n(\boldsymbol{\theta}_n^*) (\tilde{\boldsymbol{\theta}}_n - \hat{\boldsymbol{\theta}}_n), \quad (27)$$

where $\boldsymbol{\theta}_n^*$ is between $\tilde{\boldsymbol{\theta}}_n$ and $\hat{\boldsymbol{\theta}}_n$. Thus $\boldsymbol{\theta}_n^* \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$ and Corollary 1 implies $I_n(\boldsymbol{\theta}_n^*) \xrightarrow[n \rightarrow \infty]{P} I(\boldsymbol{\theta}_X)$.

Further by Theorem 5

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = I^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X) + o_P(1),$$

which together with (25) implies

$$\begin{aligned} \sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \widetilde{\boldsymbol{\theta}}_n) &= \sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) + \sqrt{n}(\boldsymbol{\theta}_X - \widetilde{\boldsymbol{\theta}}_n) \\ &= I^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X) - \left(\begin{array}{cc} \mathbf{0}_q & \\ I_{22}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_{2n}(\boldsymbol{\theta}_X) & \end{array} \right) + o_P(1) \\ &= \mathbb{A}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X) + o_P(1), \end{aligned}$$

where

$$\mathbb{A}(\boldsymbol{\theta}_X) = I^{-1}(\boldsymbol{\theta}_X) - \left(\begin{array}{cc} \mathbf{0}_{q \times q} & \mathbf{0}_{q \times (p-q)} \\ \mathbf{0}_{(p-q) \times q} & I_{22}^{-1}(\boldsymbol{\theta}_X) \end{array} \right).$$

By the central limit theorem (for i.i.d. vectors) and the symmetry of matrix $\mathbb{A}(\boldsymbol{\theta}_X)$

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \widetilde{\boldsymbol{\theta}}_n) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{A}(\boldsymbol{\theta}_X) I(\boldsymbol{\theta}_X) \mathbb{A}(\boldsymbol{\theta}_X)). \quad (28)$$

Now we will use Lemma A6 about the distribution of a quadratic form from Appendix.

Put

$$\mathbb{B} = I(\boldsymbol{\theta}_X) \quad \text{and} \quad \mathbb{V} = \mathbb{A}(\boldsymbol{\theta}_X) I(\boldsymbol{\theta}_X) \mathbb{A}(\boldsymbol{\theta}_X).$$

Now $\mathbb{B}\mathbb{V} = I(\boldsymbol{\theta}_X) \mathbb{A}(\boldsymbol{\theta}_X) I(\boldsymbol{\theta}_X) \mathbb{A}(\boldsymbol{\theta}_X)$, where

$$\begin{aligned} I(\boldsymbol{\theta}_X) \mathbb{A}(\boldsymbol{\theta}_X) &= \begin{pmatrix} I_{11}(\boldsymbol{\theta}_X) & I_{12}(\boldsymbol{\theta}_X) \\ I_{21}(\boldsymbol{\theta}_X) & I_{22}(\boldsymbol{\theta}_X) \end{pmatrix} \left(I^{-1}(\boldsymbol{\theta}_X) - \begin{pmatrix} \mathbf{0}_{q \times q} & \mathbf{0}_{q \times (p-q)} \\ \mathbf{0}_{(p-q) \times q} & I_{22}^{-1}(\boldsymbol{\theta}_X) \end{pmatrix} \right) \\ &= \mathbb{I}_p - \underbrace{\begin{pmatrix} \mathbf{0}_{q \times q} & I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) \\ \mathbf{0}_{(p-q) \times q} & \mathbb{I}_{p-q} \end{pmatrix}}_{=: \mathbb{D}}. \end{aligned}$$

Note that matrix \mathbb{D} is idempotent, thus also $\mathbb{I}_p - \mathbb{D}$ and $\mathbb{B}\mathbb{V} = (\mathbb{I}_p - \mathbb{D})(\mathbb{I}_p - \mathbb{D})$ are idempotent.

Now using (27), (28), CS (Theorem 2), Lemma A6 and CMT (Theorem 1) one gets

$$LR_n^* = 2 \left(\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\widetilde{\boldsymbol{\theta}}_n) \right) = \sqrt{n}(\widetilde{\boldsymbol{\theta}}_n - \widehat{\boldsymbol{\theta}}_n)^\top I(\boldsymbol{\theta}_X) \sqrt{n}(\widetilde{\boldsymbol{\theta}}_n - \widehat{\boldsymbol{\theta}}_n) + o_P(1) \xrightarrow[n \rightarrow \infty]{d} \chi_{\text{tr}(\mathbb{B}\mathbb{V})}^2,$$

where $\text{tr}(\mathbb{B}\mathbb{V}) = \text{tr}(\mathbb{I}_p) - \text{tr}(\mathbb{D}) = p - (p - q) = q$. \square

Suppose that both $\widehat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta})$ and $\widetilde{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta_0} \ell_n(\boldsymbol{\theta})$ (where Θ_0 stands for the parameter space under the null hypothesis) are consistent estimator under the null hypothesis. Then the likelihood ratio test can be rewritten as

$$LR_n^* = 2 \left(\ell_n(\widehat{\boldsymbol{\theta}}_n) - \ell_n(\widetilde{\boldsymbol{\theta}}_n) \right) = 2 \left(\sup_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}) - \sup_{\boldsymbol{\theta} \in \Theta_0} \ell_n(\boldsymbol{\theta}) \right). \quad (29)$$

So with the likelihood ratio test one does not need to bother with the parametrization of the parametric spaces Θ and Θ_0 so that it fits into the framework of testing $H_0 : \boldsymbol{\tau}_X = \boldsymbol{\tau}_0$. The degrees of freedom of the asymptotic distribution are determined as the difference of the dimensions of the parametric spaces Θ and Θ_0 .

Example 23. The following data gives the number of male children among the first 12 children of family size 13 in 6115 families taken from hospital records in the 19th century Saxony. The 13th child is ignored to assuage the effect of families non-randomly stopping when a desired gender is reached. Test the null hypothesis that the gender of the babies can

Nr. of boys	0	1	2	3	4	5	6	7	8	9	10	11	12
Nr. of families	3	24	104	286	670	1033	1343	1112	829	478	181	45	7

be viewed as realisations of independent random variables having the same probability of a baby boy for each family.

Hint. Let X_i stand for the number of boys in the i -th family ($i \in \{1, \dots, n\}$, where n stands for the sample size). Then the counts in the table can be represented by

$$n_k = \sum_{i=1}^n \mathbb{1}\{X_i = k\}, \quad k \in \{0, 1, \dots, 12\}$$

and the table can be viewed as a realisation of a random vector $(n_0, n_1, \dots, n_{12})^\top$ that follows multinomial distribution $\text{Mult}_{13}(n, \boldsymbol{\pi})$.

Note that under the null hypothesis X_i follows the binomial distribution, thus

$$\pi_k = \mathbb{P}(X_i = k) = \binom{12}{k} p^k (1-p)^{12-k}, \quad k \in \{0, 1, \dots, 12\},$$

where $p \in (0, 1)$ is the probability of baby boy.

Thus to parametrize the problem (so that it fits into the framework of this section) put $\psi = p$ and get

$$\pi_0 = (1 - \psi)^{12}, \quad \pi_k = \binom{12}{k} \psi^k (1 - \psi)^{12-k} + \tau_k, \quad k \in \{1, \dots, 11\},$$

and $\pi_{12} = 1 - \sum_{k=0}^{11} \pi_k$. The hypotheses can now be written as

$$H_0 : (\tau_1, \dots, \tau_{11})^\top = \mathbf{0}_{11}, \quad H_1 : (\tau_1, \dots, \tau_{11})^\top \neq \mathbf{0}_{11}.$$

Nevertheless it would be rather tedious to derive either the Wald statistic (W_n^*) or Rao score statistic (R_n^*) as one needs to calculate the score statistic and (empirical) Fisher information matrix.

On the other hand using (29) it is straightforward to calculate the likelihood ratio test LR_n^* as

$$\sup_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}) = \sum_{k=0}^{12} n_k \log \left(\frac{n_k}{n} \right)$$

and

$$\sup_{\boldsymbol{\theta} \in \Theta_0} \ell_n(\boldsymbol{\theta}) = \sum_{k=0}^{12} n_k \log \tilde{\pi}_k, \quad \text{where} \quad \tilde{\pi}_k = \binom{12}{k} (\tilde{\psi}_n)^k (1 - \tilde{\psi}_n)^{12-k}, \quad \text{with} \quad \tilde{\psi}_n = \sum_{k=1}^{12} \frac{k n_k}{12n}.$$

By Theorem 7 under the null hypothesis the test statistic LR_n^* converges in distribution to χ^2 -distribution with 11 degrees of freedom.

Another approach to test the hypothesis of interest would be (to forget about the test statistics LR_n^* , W_n^* , R_n^* and) to use the standard χ^2 -test of goodness-of-fit in multinomial distribution with estimated parameters. The test statistics would be

$$X^2 = \sum_{k=0}^{12} \frac{(n_k - n \tilde{\pi}_k)^2}{n \tilde{\pi}_k} \quad (30)$$

and under the null hypothesis it has also asymptotically χ^2 -distribution with 11 degrees of freedom. In fact it can be proved* that the test statistic X^2 given by (30) corresponds to the test statistic of the Rao score test (R_n^*) with $I^{11}(\tilde{\boldsymbol{\theta}}_n)$ taken as \hat{I}_n^{11} .

Example 24. Breusch-Pagan test of heteroscedasticity.

Example 25. Suppose that you observe independent identically distributed random vectors $\begin{pmatrix} Y_1 \\ \mathbf{X}_1 \end{pmatrix}, \dots, \begin{pmatrix} Y_n \\ \mathbf{X}_n \end{pmatrix}$ such that

$$P(Y_1 = 1 | \mathbf{X}_1) = \frac{\exp\{\alpha + \mathbf{X}_1^\top \boldsymbol{\beta}\}}{1 + \exp\{\alpha + \mathbf{X}_1^\top \boldsymbol{\beta}\}}, \quad P(Y_1 = 0 | \mathbf{X}_1) = \frac{1}{1 + \exp\{\alpha + \mathbf{X}_1^\top \boldsymbol{\beta}\}},$$

where the distribution of $\mathbf{X}_1 = (X_{11}, \dots, X_{1p})^\top$ does not depend on the unknown parameters α a $\boldsymbol{\beta}$.

- (i) Derive a test for the null hypothesis $H_0 : \boldsymbol{\beta} = \mathbf{0}_p$ against the alternative that $H_1 : \boldsymbol{\beta} \neq \mathbf{0}_p$.
- (ii) Find the confidence set for the parameter $\boldsymbol{\beta}$.

Literature: Anděl [2007] Chapter 8.6, Kulich [2014], Zvára [2008] pp. 122–128.

The end of
class 7
(12.3.2024)

* More precisely, it is said so in the textbooks but I have not managed to find the derivation.

2.7 Profile likelihood*

Let $\boldsymbol{\theta}$ be divided into $\boldsymbol{\tau}$ containing the first q components ($1 \leq q < p$) and $\boldsymbol{\psi}$ containing the remaining $p - q$ components, i.e.

$$\boldsymbol{\theta} = (\boldsymbol{\tau}^\top, \boldsymbol{\psi}^\top)^\top = (\theta_1, \dots, \theta_q, \theta_{q+1}, \dots, \theta_p)^\top.$$

Write the likelihood of the parameter $\boldsymbol{\theta}$ as $L_n(\boldsymbol{\theta}) = L_n(\boldsymbol{\tau}, \boldsymbol{\psi})$ and analogously for log-likelihood, score function, Fisher information matrix, ...

In this subsection we will assume that there exists $\hat{\boldsymbol{\theta}}_n$ which is a unique maximum of the function $\ell_n(\boldsymbol{\theta})$ and also a consistent estimator of $\boldsymbol{\theta}_X$. Similarly for $\boldsymbol{\tau}_X = \boldsymbol{\tau}$ let $\tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})$ be a unique maximum of $\ell_n(\boldsymbol{\tau}, \boldsymbol{\psi})$ and a consistent estimator of $\boldsymbol{\psi}_X$.

The profile likelihood and the profile log-likelihood for the parameter $\boldsymbol{\tau}$ are defined subsequently as

$$L_n^{(p)}(\boldsymbol{\tau}) = \max_{\boldsymbol{\psi} \in \Theta_\psi} L_n(\boldsymbol{\tau}, \boldsymbol{\psi}), \quad \ell_n^{(p)}(\boldsymbol{\tau}) = \log L_n^{(p)}(\boldsymbol{\tau}) = \max_{\boldsymbol{\psi} \in \Theta_\psi} \ell_n(\boldsymbol{\tau}, \boldsymbol{\psi}).$$

In the following we will show that one can work with the profile likelihood as with the 'standard' likelihood.

First of all put

$$\hat{\boldsymbol{\tau}}_n^{(p)} = \arg \max_{\boldsymbol{\tau} \in \Theta_\tau} \ell_n^{(p)}(\boldsymbol{\tau}).$$

Note that

$$\ell_n^{(p)}(\hat{\boldsymbol{\tau}}_n^{(p)}) = \max_{\boldsymbol{\tau} \in \Theta_\tau} \ell_n^{(p)}(\boldsymbol{\tau}) = \max_{\boldsymbol{\tau} \in \Theta_\tau} \max_{\boldsymbol{\psi} \in \Theta_\psi} \ell_n(\boldsymbol{\tau}, \boldsymbol{\psi}) = \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}) = \ell_n(\hat{\boldsymbol{\theta}}_n).$$

As we assume that $\hat{\boldsymbol{\theta}}_n$ is a unique maximizer of $\ell_n(\boldsymbol{\theta})$, this implies that

$$\hat{\boldsymbol{\tau}}_n^{(p)} = \hat{\boldsymbol{\tau}}_n,$$

where $\hat{\boldsymbol{\tau}}_n$ stands for the first q -coordinates of the maximum likelihood estimator $\hat{\boldsymbol{\theta}}_n$.

Further denote

$$\tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}) = \arg \max_{\boldsymbol{\psi} \in \Theta_\psi} \ell_n(\boldsymbol{\tau}, \boldsymbol{\psi}), \quad \tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau}) = (\boldsymbol{\tau}^\top, \tilde{\boldsymbol{\psi}}_n^\top(\boldsymbol{\tau}))^\top,$$

and define the profile score statistic and profile (empirical) information matrix as

$$\mathbf{U}_n^{(p)}(\boldsymbol{\tau}) = \frac{\partial \ell_n^{(p)}(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}}, \quad I_n^{(p)}(\boldsymbol{\tau}) = -\frac{1}{n} \frac{\partial \mathbf{U}_n^{(p)}(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top}.$$

The following lemma shows how the quantities $\mathbf{U}_n^{(p)}(\boldsymbol{\tau})$ and $I_n^{(p)}(\boldsymbol{\tau})$ are related with $\mathbf{U}_n(\boldsymbol{\theta})$ and $I_n(\boldsymbol{\theta})$.

* *Profilová věrohodnost.*

Lemma 3. Suppose that assumptions **[R0]**-**[R6]** are satisfied. Then (with probability tending to one) on a neighbourhood of τ_X

$$\mathbf{U}_n^{(p)}(\boldsymbol{\tau}) = \mathbf{U}_{1n}(\tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau})), \quad I_n^{(p)}(\boldsymbol{\tau}) = I_{11n}(\tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau})) - I_{12n}(\tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau}))I_{22n}^{-1}(\tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau}))I_{21n}(\tilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau})),$$

where $I_{jkn}(\boldsymbol{\theta})$ (for $j, k \in \{1, 2\}$) were introduced in (20).

Proof. $\mathbf{U}_n^{(p)}(\boldsymbol{\tau})$: Let us calculate

$$\begin{aligned} [\mathbf{U}_n^{(p)}(\boldsymbol{\tau})]^\top &= \frac{\partial \ell_n^{(p)}(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top} = \frac{\partial \ell_n(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}))}{\partial \boldsymbol{\tau}^\top} \\ &= \mathbf{U}_{1n}^\top(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) + \mathbf{U}_{2n}^\top(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) \frac{\partial \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top} = \mathbf{U}_{1n}^\top(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})), \end{aligned} \quad (31)$$

where the last equality follows from the fact that $\tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}) = \arg \max_{\boldsymbol{\psi} \in \Theta_\psi} \ell_n^{(p)}(\boldsymbol{\tau}, \boldsymbol{\psi})$, which implies that $\mathbf{U}_{2n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) = \mathbf{0}_{p-q}$.

$I_n^{(p)}(\boldsymbol{\tau})$: Note that with the help of (31)

$$\begin{aligned} I_n^{(p)}(\boldsymbol{\tau}) &= -\frac{1}{n} \frac{\partial \mathbf{U}_n^{(p)}(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top} = -\frac{1}{n} \frac{\partial \mathbf{U}_{1n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}))}{\partial \boldsymbol{\tau}^\top} \\ &= I_{11,n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) + I_{12,n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) \frac{\partial \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top}. \end{aligned} \quad (32)$$

Further by differentiating both sides of the identity

$$\mathbf{U}_{2n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) = \mathbf{0}_{p-q}$$

with respect to $\boldsymbol{\tau}^\top$ one gets

$$I_{21,n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) + I_{22,n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) \frac{\partial \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top} = \mathbf{0}_{(p-q) \times q},$$

which implies that

$$\frac{\partial \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}^\top} = -I_{22,n}^{-1}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})) I_{21,n}(\boldsymbol{\tau}, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau})). \quad (33)$$

Now combining (32) and (33) implies the statement of the theorem for $I_n^{(p)}(\boldsymbol{\tau})$. \square

Tests based on profile likelihood

Define the (profile) test statistics of the null hypothesis $H_0 : \tau_X = \tau_0$ as

$$\begin{aligned} LR_n^{(p)} &= 2(\ell_n^{(p)}(\hat{\boldsymbol{\tau}}_n) - \ell_n^{(p)}(\boldsymbol{\tau}_0)), \\ W_n^{(p)} &= n(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0)^\top \hat{I}_n^{(p)}(\hat{\boldsymbol{\tau}}_n)(\hat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0), \\ R_n^{(p)} &= \frac{1}{n} [\mathbf{U}_n^{(p)}(\boldsymbol{\tau}_0)]^\top [\hat{I}_n^{(p)}]^{-1} \mathbf{U}_n^{(p)}(\boldsymbol{\tau}_0), \end{aligned}$$

where one can use for instance $I_n^{(p)}(\boldsymbol{\tau}_0)$ or $I_n^{(p)}(\hat{\boldsymbol{\tau}}_n)$ as $\hat{I}_n^{(p)}$.

Theorem 8. *Suppose that the null hypothesis holds and assumptions [R0]-[R6] are satisfied. Then each of the test statistics $LR_n^{(p)}$, $W_n^{(p)}$ and $R_n^{(p)}$ converges in distribution to χ^2 -distribution with q degrees of freedom.*

Proof. $LR_n^{(p)}$: Note that

$$\ell_n^{(p)}(\widehat{\boldsymbol{\tau}}_n) = \ell_n(\widehat{\boldsymbol{\tau}}_n, \widehat{\boldsymbol{\psi}}_n) = \ell_n(\widehat{\boldsymbol{\theta}}_n)$$

and further

$$\ell_n^{(p)}(\boldsymbol{\tau}_0) = \max_{\boldsymbol{\psi} \in \Theta_{\boldsymbol{\psi}}} \ell_n(\boldsymbol{\tau}_0, \boldsymbol{\psi}) = \ell_n(\boldsymbol{\tau}_0, \widetilde{\boldsymbol{\psi}}_n) = \ell_n(\widetilde{\boldsymbol{\theta}}_n).$$

Thus $LR_n^{(p)} = LR_n^*$, where LR_n^* is the test statistic of the likelihood ratio test in the presence of nuisance parameters given by (21). Thus the statement of the theorem follows by Theorem 7.

$W_n^{(p)}$: Follows from Theorem 7 and the fact that by Lemmas 1, 2 and 3

$$\widehat{I}_n^{(p)} \xrightarrow[n \rightarrow \infty]{P} I_{11}(\boldsymbol{\theta}_X) - I_{12}(\boldsymbol{\theta}_X) I_{22}^{-1}(\boldsymbol{\theta}_X) I_{21}(\boldsymbol{\theta}_X) = [I^{11}(\boldsymbol{\theta}_X)]^{-1}. \quad (34)$$

$R_n^{(p)}$: By Lemma 3 one has $\mathbf{U}_n^{(p)}(\boldsymbol{\tau}) = \mathbf{U}_{1n}(\widetilde{\boldsymbol{\theta}}_n(\boldsymbol{\tau}))$. Thus $R_n^{(p)} = R_n^*$ with $\widehat{I}_n^{11} = [\widehat{I}_n^{(p)}]^{-1}$, where R_n^* is Rao score test statistic in the presence of nuisance parameters defined in (22). The statement of the theorem now follows by (34) and Theorem 7. \square

Confidence interval for θ_{X_j}

One of the applications of the profile likelihood is to construct a confidence interval for θ_{X_j} . Let $\tau = \theta_j$ and $\boldsymbol{\psi}$ contains the remaining coordinates of the parameter $\boldsymbol{\theta}$. Then the set

$$\left\{ \theta_j : 2 \left(\ell_n^{(p)}(\widehat{\boldsymbol{\theta}}_{nj}) - \ell_n^{(p)}(\theta_j) \right) \leq \chi_1^2(1 - \alpha) \right\}$$

is the asymptotic confidence interval for θ_{X_j} . Although this confidence interval is more difficult to calculate than the Wald-type confidence interval given by (19), the simulations show that it has better finite sample properties. In R-software these intervals for GLM models are calculated by the function `confint`.

Example 26. Let X_1, \dots, X_n be a random sample from a gamma distribution with density

$$f(x) = \frac{1}{\Gamma(\beta)} \lambda^\beta x^{\beta-1} \exp\{-\lambda x\} \mathbb{1}\{x > 0\}.$$

Suppose we are interested in parameter β and parameter λ is nuisance. Derive the profile likelihood for parameter β and the Rao score test of the null hypothesis $H_0 : \beta_X = \beta_0$ against $H_1 : \beta_X \neq \beta_0$ that is based on the profile likelihood.

Solution: The likelihood and log-likelihood are given by

$$L_n(\beta, \lambda) = \prod_{i=1}^n \frac{1}{\Gamma(\beta)} \lambda^\beta X_i^{\beta-1} e^{-\lambda X_i},$$

$$\ell_n(\beta, \lambda) = -n \log \Gamma(\beta) + n\beta \log \lambda + (\beta - 1) \sum_{i=1}^n \log X_i - \lambda \sum_{i=1}^n X_i.$$

For a given β we can find $\tilde{\lambda}_n(\beta)$ by

$$\begin{aligned} \frac{\partial \ell_n(\beta, \lambda)}{\partial \lambda} &= \frac{n\beta}{\lambda} - \sum_{i=1}^n X_i \stackrel{!}{=} 0 \\ \tilde{\lambda}_n(\beta) &= \frac{\beta}{\bar{X}_n}. \end{aligned}$$

Thus the profile log-likelihood is

$$\ell_n^{(p)}(\beta) = -n \log \Gamma(\beta) + n\beta \log \left(\frac{\beta}{\bar{X}_n} \right) + (\beta - 1) \sum_{i=1}^n \log X_i - n\beta$$

and its corresponding score function

$$U_n^{(p)}(\beta) = -\frac{n \Gamma'(\beta)}{\Gamma(\beta)} + n \log \left(\frac{\beta}{\bar{X}_n} \right) + n + \sum_{i=1}^n \log X_i - n.$$

Statistic of Rao score test of the null hypothesis $H_0 : \beta_X = \beta_0$ against $H_1 : \beta_X \neq \beta_0$ is now given by

$$R_n^{(p)} = \frac{[U_n^{(p)}(\beta_0)]^2}{n I_n^{(p)}(\beta_0)},$$

where

$$I_n^{(p)}(\beta) = -\frac{1}{n} \frac{\partial U_n^{(p)}(\beta)}{\partial \beta} = \left[\frac{\Gamma''(\beta)}{\Gamma(\beta)} - \left(\frac{\Gamma'(\beta)}{\Gamma(\beta)} \right)^2 - \frac{1}{\beta} \right].$$

Example 27. Box-Cox transformation. See [Zvára \[2008\]](#) pp. 149–151.

Remark 12. Although we have shown that one can work with the profile likelihood as with the standard likelihood not all the properties are shared. For instance for standard score statistic one has $\mathbf{E} \mathbf{U}_n(\boldsymbol{\theta}_X) = \mathbf{0}_p$. But this is not guaranteed for profile score statistic as by [Lemma 3](#)

$$\mathbf{E} \mathbf{U}_n^{(p)}(\boldsymbol{\tau}_X) = \mathbf{E} \mathbf{U}_{1n}(\boldsymbol{\tau}_X, \tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}_X))$$

and the expectation on the right-hand side of the previous equation is typically not zero due to the random argument $\tilde{\boldsymbol{\psi}}_n(\boldsymbol{\tau}_X)$ (for illustration think of $\mathbf{E} U_n^{(p)}(\beta_X)$ in [Example 26](#)). From the proof of [Theorem 7](#) we only know that $\frac{1}{\sqrt{n}} \mathbf{U}_n^{(p)}(\boldsymbol{\tau}_X)$ converges in distribution to a zero-mean Gaussian distribution.

Note also that we have avoided defining the profile Fisher information matrix. The thing is that the only definition that makes sense would be $I^{(p)}(\boldsymbol{\tau}_X) = [I^{11}(\boldsymbol{\tau}_X, \boldsymbol{\psi}_X)]^{-1}$. But this is not nice as it depends on the nuisance parameter $\boldsymbol{\psi}_X$. Further, it does not hold that $I^{(p)}(\boldsymbol{\tau}_X)$ is the expectation of $I_n^{(p)}(\boldsymbol{\tau}_X)$. It only holds that

$$I_n^{(p)}(\boldsymbol{\tau}_X) \xrightarrow[n \rightarrow \infty]{P} I^{(p)}(\boldsymbol{\tau}_X).$$

2.8 Some notes on maximum likelihood in case of not i.i.d. random vectors

Let observations $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$ have a joint density $f_n(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta})$ that is known up to the unknown parameter $\boldsymbol{\theta}$ from the parametric space Θ . Analogously as in ‘i.i.d case’ one can define the *likelihood function* as

$$L_n(\boldsymbol{\theta}) = f_n(\mathbf{X}_1, \dots, \mathbf{X}_n; \boldsymbol{\theta}),$$

the *log-likelihood function* as

$$\ell_n(\boldsymbol{\theta}) = \log L_n(\boldsymbol{\theta}),$$

and the *score statistic* as

$$\mathbf{U}_n(\boldsymbol{\theta}) = \frac{\partial \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

The *maximum likelihood estimator* (of parameter $\boldsymbol{\theta}_X$) is then defined as

$$\widehat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} L_n(\boldsymbol{\theta}) \quad \text{or more generally as} \quad \mathbf{U}_n(\widehat{\boldsymbol{\theta}}_n) \stackrel{!}{=} \mathbf{0}_p.$$

Finally the *observed (empirical) Fisher information matrix* as

$$I_n(\boldsymbol{\theta}) = -\frac{1}{n} \frac{\partial^2 \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}.$$

The role of the theoretical Fisher information matrix $I(\boldsymbol{\theta})$ in ‘i.i.d’ settings is now taken by the *limit ‘average’ Fisher information matrix*

$$\bar{I}(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[\frac{-\partial^2 \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right].$$

In ‘nice (regular) models’ (see also Remark 13 below) it holds that

$$\sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}, \bar{I}^{-1}(\boldsymbol{\theta}_X)).$$

The most straightforward estimator of $\bar{I}(\boldsymbol{\theta}_X)$ is $I_n(\widehat{\boldsymbol{\theta}}_n)$ and thus the estimator of the asymptotic variance matrix of $\widehat{\boldsymbol{\theta}}_n$ is

$$\widehat{\text{avar}}(\widehat{\boldsymbol{\theta}}_n) = \frac{1}{n} I_n^{-1}(\widehat{\boldsymbol{\theta}}_n) = \left[\frac{-\partial^2 \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_n} \right]^{-1}.$$

That is why some authors prefer to define the empirical Fisher information without $\frac{1}{n}$ simply as

$$\tilde{I}_n(\boldsymbol{\theta}) = \frac{-\partial^2 \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

and they speak about it as the Fisher information of all observations.

Remark 13. An inspection of the proof of Theorem 9 (for Z -estimators) reveals that we need to show the analogy of Lemma 1 with $I(\boldsymbol{\theta}_X)$ replaced with $\bar{I}(\boldsymbol{\theta}_X)$ and that

$$\frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \bar{I}(\boldsymbol{\theta}_X)).$$

Example 28. Suppose we have K independent samples, that is for each $k \in \{1, \dots, K\}$ the random variables $\mathbf{X}_{ki}, i \in \{1, \dots, n_k\}$ are independent and identically distributed with density $f_k(\mathbf{x}; \boldsymbol{\theta})$ (with respect to a σ -finite measure μ). Further let all the random variables be independent and let $\lim_{n \rightarrow \infty} \frac{n_k}{n} = w_k$, where $n = n_1 + \dots + n_K$. Then

$$\begin{aligned} L_n(\boldsymbol{\theta}) &= \prod_{k=1}^K \prod_{i=1}^{n_k} f_k(\mathbf{X}_{ki}; \boldsymbol{\theta}), \\ \ell_n(\boldsymbol{\theta}) &= \sum_{k=1}^K \sum_{i=1}^{n_k} \log f_k(\mathbf{X}_{ki}; \boldsymbol{\theta}), \\ \mathbf{U}_n(\boldsymbol{\theta}) &= \frac{\partial \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{k=1}^K \sum_{i=1}^{n_k} \frac{\partial \log f_k(\mathbf{X}_{ki}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \\ I_n(\boldsymbol{\theta}) &= -\frac{1}{n} \frac{\partial \mathbf{U}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} = -\frac{1}{n} \sum_{k=1}^K \sum_{i=1}^{n_k} \frac{\partial^2 \log f_k(\mathbf{X}_{ki}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}, \\ \bar{I}(\boldsymbol{\theta}) &= \lim_{n \rightarrow \infty} \mathbb{E} I_n(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} \sum_{k=1}^K \underbrace{\frac{n_k}{n}}_{\rightarrow w_k} I^{(k)}(\boldsymbol{\theta}) = \sum_{k=1}^K w_k I^{(k)}(\boldsymbol{\theta}), \end{aligned}$$

where $I^{(k)}(\boldsymbol{\theta})$ is Fisher information matrix of \mathbf{X}_{ki} (i.e. for the density $f_k(\mathbf{x}; \boldsymbol{\theta})$).

In standard applications $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_K^\top)^\top$, and the density $f_k(\mathbf{x}; \boldsymbol{\theta})$ depends only on $\boldsymbol{\theta}_k$, i.e. $f_k(\mathbf{x}; \boldsymbol{\theta}) = f(\mathbf{x}; \boldsymbol{\theta}_k)$. And we are usually interested in testing the null hypothesis that all the distributions are the same, that is

$$H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 = \dots = \boldsymbol{\theta}_K \quad H_1 : \exists_{k,j \in \{1, \dots, K\}} \boldsymbol{\theta}_k \neq \boldsymbol{\theta}_j.$$

See also Example 32.

Random vs. fixed design

Sometimes in regression it is useful to distinguish random design and fixed design.

In **random design** we assume that the values of the covariates are realisations of random vectors. Thus (in the most simple situation) we assume that we observe independent and identically distributed random vectors

$$\begin{pmatrix} Y_1 \\ \mathbf{X}_1 \end{pmatrix}, \dots, \begin{pmatrix} Y_n \\ \mathbf{X}_n \end{pmatrix}, \quad (35)$$

where the conditional distribution of $Y_i|\mathbf{X}_i$ is known up to the unknown parameter $\boldsymbol{\theta}$ and the distribution of \mathbf{X}_i does not depend on $\boldsymbol{\theta}$. Put $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ for the conditional density of $Y_i|\mathbf{X}_i = \mathbf{x}_i$ and $f_{\mathbf{X}}(\mathbf{x})$ for the density of \mathbf{X}_i . Then the likelihood and the log-likelihood (for the parameter $\boldsymbol{\theta}$) are given by

$$\begin{aligned} L_n(\boldsymbol{\theta}) &= \prod_{i=1}^n f_{Y,\mathbf{X}}(Y_i, \mathbf{X}_i; \boldsymbol{\theta}) = \prod_{i=1}^n f(Y_i|\mathbf{X}_i; \boldsymbol{\theta}) f_{\mathbf{X}}(\mathbf{X}_i) \\ \ell_n(\boldsymbol{\theta}) &= \sum_{i=1}^n \log f(Y_i|\mathbf{X}_i; \boldsymbol{\theta}) + \sum_{i=1}^n \log f_{\mathbf{X}}(\mathbf{X}_i). \end{aligned} \quad (36)$$

In **fixed design** it is assumed that the values of the covariates $\mathbf{x}_1, \dots, \mathbf{x}_n$ are fixed when planning the experiment (before measuring the response). Now we observe Y_1, \dots, Y_n independent (but not identically distributed) random variables with the densities $f(y_1|\mathbf{x}_1; \boldsymbol{\theta}), \dots, f(y_n|\mathbf{x}_n; \boldsymbol{\theta})$. Then the log-likelihood is given by

$$\ell_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(Y_i|\mathbf{x}_i; \boldsymbol{\theta}). \quad (37)$$

Comparing the log-likelihoods in (36) and (37) one can see that (once the data are observed) they differ only by $\sum_{i=1}^n \log f_{\mathbf{X}}(\mathbf{X}_i)$ which does not depend on $\boldsymbol{\theta}$. Thus in terms of (likelihood based) inference for a given dataset both approaches are equivalent. The only difference is that the theory for the fixed design is more difficult.

Example 29. *Poisson regression.*

Random design approach: We assume that we observe independent identically distributed random vectors (35) and that $Y_i|\mathbf{X}_i \sim \text{Po}(\lambda(\mathbf{X}_i))$, where $\lambda(\mathbf{x}) = \exp\{\mathbf{x}^\top \boldsymbol{\beta}\}$ and $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$. Then (provided assumptions **[R0]**-**[R6]** are satisfied)

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, I^{-1}(\boldsymbol{\beta}_X)), \text{ where } I(\boldsymbol{\beta}_X) = \mathbf{E} [\mathbf{X}_1 \mathbf{X}_1^\top \exp\{\mathbf{X}_1^\top \boldsymbol{\beta}_X\}].$$

Fixed design approach: We assume that we observe independent random variables Y_1, \dots, Y_n and we have the known constants $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that $Y_i \sim \text{Po}(\lambda(\mathbf{x}_i))$, where $\lambda(\mathbf{x}) =$

$\exp\{\mathbf{x}^\top \boldsymbol{\beta}\}$. Then it can be shown (that under mild assumptions on $\mathbf{x}_1, \dots, \mathbf{x}_n$)

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \bar{I}^{-1}(\boldsymbol{\beta}_X)), \text{ where } \bar{I}(\boldsymbol{\beta}_X) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}_X\}.$$

Note that in practice both $I(\boldsymbol{\beta}_X)$ and $\bar{I}(\boldsymbol{\beta}_X)$ would be estimated by

$$\widehat{I}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \exp\{\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n\} \quad \text{or} \quad \widehat{\bar{I}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \exp\{\mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}_n\}.$$

Thus for observed data the estimators coincide. The only difference is in notation in which you distinguish whether you think of the observed values of the covariates as the realizations of the random vectors or as fixed constants.

Example 30. Note that alternatively one can view the K-sample problem described in Example 28 also within i.i.d framework. Consider the data as a realization of the random sample $(\mathbf{Z}_{J_1}^1), \dots, (\mathbf{Z}_{J_n}^n)$, where J_i takes values in $\{1, \dots, K\}$ and the conditional distribution of \mathbf{Z}_i given $J_i = j$ is given by the density $f_j(\mathbf{x}; \boldsymbol{\theta})$.

Example 31. Maximum likelihood estimation in AR(1) process.

Example 32. Suppose that X_{ki} , $k \in \{1, \dots, K\}$, $i = 1, \dots, n_K$ be independent random variables such that X_{ki} follows Bernoulli distribution with parameter p_k . We are interested in testing the hypothesis

$$H_0 : p_1 = p_2 = \dots = p_K \quad H_1 : \exists_{k,j \in \{1, \dots, K\}} p_k \neq p_j.$$

Note that one can easily construct a likelihood ratio test.

Alternatively one can view the data as $K \times 2$ contingency table and use the χ^2 -test of independence. It can be proved that this test is in fact the Rao-score test for this problem.

Literature: Hoadley [1971].

The end of
class 9
(19. 3. 2024)

2.9 Conditional and marginal likelihood*

In some models the number of parameters is increasing as the sample size increases. Formally let $\boldsymbol{\theta}^{(n)} = (\theta_1, \dots, \theta_{p_n})^\top$, where p_n is a non-decreasing function of n . Let $\boldsymbol{\theta}^{(n)}$ be divided into $\boldsymbol{\tau}$ containing the first q components (with q being fixed) and $\boldsymbol{\psi}^{(n)}$ containing the remaining $(p_n - q)$ components.

Example 33. *Strongly stratified sample.* Let Y_{ij} , $i \in \{1, \dots, N\}$, $j \in \{1, 2\}$ be independent random variables such that $Y_{ij} \sim \mathbf{N}(\mu_i, \sigma^2)$. Derive the maximum likelihood estimator of σ^2 . Is this estimator consistent as $N \rightarrow \infty$?

* *Podmíněná a marginální věrohodnost.*

Solution. The joint density of all the observations $\mathbf{Y} = (Y_{ij}, i \in \{1, \dots, N\}, j \in \{1, 2\})$ is

$$f(\mathbf{y}; \sigma^2, \mu_1, \dots, \mu_N) = \prod_{i=1}^N \prod_{j=1}^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(Y_{ij}-\mu_i)^2}{2\sigma^2}\right\} \quad (38)$$

and thus the log-likelihood is given by

$$\ell_n(\sigma^2, \mu_1, \dots, \mu_N) = -N \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{j=1}^2 (Y_{ij} - \mu_i)^2 - N \log(2\pi).$$

Differentiating with respect to μ_1, \dots, μ_N and σ^2 one easily finds that

$$\hat{\mu}_i = \frac{Y_{i1} + Y_{i2}}{2}, \quad i \in \{1, \dots, N\}$$

and

$$\begin{aligned} \hat{\sigma}_N^2 &= \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^2 (Y_{ij} - \hat{\mu}_i)^2 = \frac{1}{2N} \sum_{i=1}^N \left[(Y_{i1} - \frac{Y_{i1}+Y_{i2}}{2})^2 + (Y_{i2} - \frac{Y_{i1}+Y_{i2}}{2})^2 \right] \\ &= \frac{1}{2N} \sum_{i=1}^N \left[(\frac{Y_{i1}-Y_{i2}}{2})^2 + (\frac{Y_{i2}-Y_{i1}}{2})^2 \right] = \frac{1}{4N} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2. \end{aligned}$$

Thus

$$\hat{\sigma}_N^2 \xrightarrow[N \rightarrow \infty]{P} \frac{1}{4} \mathbf{E} (Y_{i1} - Y_{i2})^2 = \frac{1}{4} \text{var}(Y_{i1} - Y_{i2}) = \frac{2\sigma^2}{4} = \frac{\sigma^2}{2} \neq \sigma^2.$$

Note that in Example 33 each observation carries information on σ^2 , but the maximum likelihood estimator of σ^2 is not even consistent. The problem is that the dimension of nuisance parameters $\boldsymbol{\psi}^{(N)} = (\mu_1, \dots, \mu_N)^\top$ is increasing to infinity (too quickly). Marginal and conditional likelihoods are two attempts to modify the likelihood so that it yields a consistent (and hopefully also asymptotically normal) estimator of the parameter $\boldsymbol{\tau}$.

Suppose that one can use data \mathbb{X} to calculate \mathbf{V} whose distribution depends only on parameter $\boldsymbol{\tau}$ (and not on $\boldsymbol{\psi}^{(n)}$). Then *the marginal (log-)likelihood* of parameter $\boldsymbol{\tau}$ is defined as

$$L_n^{(M)}(\boldsymbol{\tau}) = f(\mathbf{V}; \boldsymbol{\tau}), \quad \ell_n^{(M)}(\boldsymbol{\tau}) = \log(L_n^{(M)}(\boldsymbol{\tau})),$$

where $f(\mathbf{v}; \boldsymbol{\tau})$ is the joint density of \mathbf{V} with respect to a σ -finite measure μ .

Suppose that one can use data \mathbb{X} to calculate \mathbf{V} and \mathbf{W} such that the conditional distribution of \mathbf{V} given \mathbf{W} depends only on parameter $\boldsymbol{\tau}$ (and not on $\boldsymbol{\psi}^{(n)}$). Then *the conditional (log-)likelihood* of parameter $\boldsymbol{\tau}$ is defined as

$$L_n^{(C)}(\boldsymbol{\tau}) = f(\mathbf{V} | \mathbf{W}; \boldsymbol{\tau}), \quad \ell_n^{(C)}(\boldsymbol{\tau}) = \log(L_n^{(C)}(\boldsymbol{\tau})),$$

where $f(\mathbf{v} | \mathbf{w}; \boldsymbol{\tau})$ is the conditional density of \mathbf{V} given $\mathbf{W} = \mathbf{w}$ with respect to a σ -finite measure μ .

Remark 14. (i) If \mathbf{V} is independent of \mathbf{W} , then $f(\mathbf{V}|\mathbf{W}; \boldsymbol{\tau}) = f(\mathbf{V}; \boldsymbol{\tau})$ and thus $L_n^{(M)}(\boldsymbol{\tau}) = L_n^{(C)}(\boldsymbol{\tau})$.

(ii) ‘Automatic calculation of $\ell_n^{(C)}(\boldsymbol{\tau})$ ’:

$$\ell_n^{(C)}(\boldsymbol{\tau}) = \log \left(\frac{f(\mathbf{V}, \mathbf{W}; \boldsymbol{\tau}, \boldsymbol{\psi}^{(n)})}{f(\mathbf{W}; \boldsymbol{\tau}, \boldsymbol{\psi}^{(n)})} \right) = \ell_{n, \mathbf{V}, \mathbf{W}}(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)}) - \ell_{n, \mathbf{W}}(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)}),$$

where $\ell_{n, \mathbf{V}, \mathbf{W}}(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)})$ is the log-likelihood of (\mathbf{V}, \mathbf{W}) and $\ell_{n, \mathbf{W}}(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)})$ is the log-likelihood of \mathbf{W} . Note that using this approach we do not need to derive the conditional distribution of \mathbf{V} given \mathbf{W} .

(iii) It can be shown that (under certain regularity assumptions) one can work with $L_n^{(M)}(\boldsymbol{\tau})$ and $L_n^{(C)}(\boldsymbol{\tau})$ as with ‘standard’ likelihoods.

The question of interest is how to find \mathbf{V} and \mathbf{W} so that we do not lose too much information about $\boldsymbol{\tau}$. To the best of my knowledge for marginal likelihood there are only ad-hoc approaches.

For conditional likelihood one can use the theory of sufficient statistics. Suppose that for each fixed value of $\boldsymbol{\tau}$ the statistic $\mathbf{S}_n(\mathbb{X})$ is sufficient for $\boldsymbol{\psi}^{(n)}$. Thus the conditional distribution of \mathbb{X} given $\mathbf{S}_n(\mathbb{X})$ does not depend on $\boldsymbol{\psi}^{(n)}$. This implies that when constructing the conditional likelihood $L_n^{(C)}(\boldsymbol{\tau})$ one can take $\mathbf{S}_n(\mathbb{X})$ as \mathbf{W} and \mathbb{X} as \mathbf{V} .

Exponential family

Let the dataset \mathbb{X} have the density (with respect to a σ -finite measure μ) of the form

$$f(\mathbf{x}; \boldsymbol{\tau}, \boldsymbol{\psi}^{(n)}) = \exp \left\{ \sum_{j=1}^q Q_j(\boldsymbol{\tau}) T_j(\mathbf{x}) + \sum_{j=1}^{p_n-q} R_j(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)}) S_j(\mathbf{x}) \right\} a(\boldsymbol{\tau}, \boldsymbol{\psi}^{(n)}) h(\mathbf{x}), \quad (39)$$

where $\boldsymbol{\tau} = (\tau_1, \dots, \tau_q)^\top$ and $\boldsymbol{\psi}^{(n)} = (\psi_1^{(n)}, \dots, \psi_{p_n-q}^{(n)})^\top$. Put $\mathbf{S}_n(\mathbb{X}) = (S_1(\mathbb{X}), \dots, S_{p_n-q}(\mathbb{X}))^\top$ and note that for a fixed value of $\boldsymbol{\tau}$ the statistic $\mathbf{S}_n(\mathbb{X})$ is sufficient for $\boldsymbol{\psi}^{(n)}$. Thus one can put $\mathbf{W} = \mathbf{S}_n(\mathbb{X})$ and $\mathbf{V} = \mathbb{X}$.

Example 33. *Strongly stratified sample (cont.).* Derive the marginal and conditional likelihood.

Marginal likelihood. For $i \in \{1, \dots, N\}$ consider $V_i = \frac{Y_{i1} - Y_{i2}}{\sqrt{2}}$. Then $V_i \sim \mathcal{N}(0, \sigma^2)$. Thus the marginal likelihood is the likelihood of V_1, \dots, V_n and is given by

$$L_n^{(M)}(\sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{V_i^2}{2\sigma^2} \right\}.$$

Further the marginal log-likelihood is given by

$$\ell_n^{(M)}(\sigma^2) = -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N V_i^2 - \frac{N}{2} \log 2\pi = -\frac{N}{2} \log \sigma^2 - \frac{1}{4\sigma^2} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2 - \frac{N}{2} \log 2\pi.$$

With this marginal log-likelihood one can work in the ‘standard’ way. That is one can for instance derive the maximum (marginal) likelihood estimator

$$\hat{\sigma}_N^{2(M)} = \frac{1}{2N} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2.$$

It is straightforward to show that this estimator is consistent and that

$$\sqrt{n} (\hat{\sigma}_N^{2(M)} - \sigma^2) \xrightarrow[N \rightarrow \infty]{d} \mathbf{N}(0, 2\sigma^4),$$

where the asymptotic variance $2\sigma^2$ can be calculated as $\text{var}((Y_{i1} - Y_{i2})^2)$ or as one over the Fisher information that corresponds to $\ell_N^{(M)}(\sigma^2)$.

Conditional likelihood. Note that the joint density (38) of $\mathbf{Y} = (Y_{ij}, i \in \{1, \dots, N\}, j \in \{1, 2\})$ can be written as

$$f(\mathbf{y}; \sigma^2, \dots) = \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{j=1}^2 y_{ij}^2 + \sum_{i=1}^N \frac{\mu_i}{\sigma^2} (y_{i1} + y_{i2}) - \sum_{i=1}^N \frac{\mu_i^2}{\sigma^2} \right\} \frac{1}{(2\pi\sigma^2)^N}. \quad (40)$$

Now the above density can be written in the form (39) with $\tau = -\frac{1}{\sigma^2}$, $\psi_i = \mu_i$, $T(\mathbf{y}) = \sum_{i=1}^N \sum_{j=1}^2 y_{ij}^2$ and $S_i(\mathbf{y}) = y_{i1} + y_{i2}$. Thus by Remark 14(ii) the conditional log-likelihood of \mathbf{Y} given $\mathbf{S}(\mathbf{Y}) = (S_1(\mathbf{Y}), \dots, S_N(\mathbf{Y}))^\top$ can be calculated as

$$\ell_n^{(C)}(\sigma^2) = \ell_n(\sigma^2, \mu_1, \dots, \mu_N) - \ell_{n, \mathbf{S}(\mathbf{Y})}(\sigma^2, \mu_1, \dots, \mu_N).$$

Here $\ell_n(\sigma^2, \mu_1, \dots, \mu_N)$ is the (standard) log-likelihood of \mathbf{Y} which can be with the help of (40) rewritten as

$$\ell_n(\sigma^2, \dots) = -N \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (Y_{i1}^2 + Y_{i2}^2) + \sum_{i=1}^N \frac{\mu_i}{\sigma^2} (Y_{i1} + Y_{i2}) - \sum_{i=1}^N \frac{\mu_i^2}{\sigma^2} - N \log 2\pi \quad (41)$$

and $\ell_{n, \mathbf{S}(\mathbf{Y})}(\sigma^2, \mu_1, \dots, \mu_N)$ is the log-likelihood of $\mathbf{S}(\mathbf{Y}) = (Y_{11} + Y_{12}, \dots, Y_{N1} + Y_{N2})^\top$. As the components of $\mathbf{S}(\mathbf{Y})$ are independent random variables with the distribution $\mathbf{N}(2\mu_i, 2\sigma^2)$ ($i \in \{1, \dots, N\}$), the log-likelihood $\ell_{n, \mathbf{S}(\mathbf{Y})}(\sigma^2, \dots)$ is given by

$$\begin{aligned} \ell_{n, \mathbf{S}(\mathbf{Y})}(\sigma^2, \dots) &= \log \left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi 2\sigma^2}} \exp \left\{ -\frac{(Y_{i1} + Y_{i2} - 2\mu_i)^2}{2 \cdot 2\sigma^2} \right\} \right) \\ &= -\frac{N}{2} \log \sigma^2 - \frac{1}{4\sigma^2} \sum_{i=1}^N (Y_{i1}^2 + Y_{i2}^2 + 2Y_{i1}Y_{i2}) \\ &\quad + \sum_{i=1}^N \frac{4\mu_i}{4\sigma^2} (Y_{i1} + Y_{i2}) - \sum_{i=1}^N \frac{4\mu_i^2}{4\sigma^2} + \frac{N}{2} \log(4\pi). \end{aligned} \quad (42)$$

Thus comparing (41) and (42) one gets

$$\begin{aligned}\ell_n^{(C)}(\sigma^2) &= \ell_n(\sigma^2, \dots) - \ell_{n, \mathbf{S}(\mathbf{Y})}(\sigma^2, \dots) = -\frac{N}{2} \log \sigma^2 - \frac{1}{4\sigma^2} \sum_{i=1}^N (Y_{i1}^2 + Y_{i2}^2 - 2Y_{i1}Y_{i2}) - \frac{N}{4} \log \pi \\ &= -\frac{N}{2} \log \sigma^2 - \frac{1}{4\sigma^2} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2 - \frac{N}{4} \log \pi = \ell_N^{(M)}(\sigma^2) + \frac{N}{4} \log \pi,\end{aligned}$$

where the (irrelevant) difference between the $\ell_n^{(C)}(\sigma^2)$ and $\ell_N^{(M)}(\sigma^2)$ comes from the fact that for the conditional likelihood we use the conditional distribution of \mathbf{Y} given $\mathbf{W} = \mathbf{S}(\mathbf{Y})$ instead of the conditional distribution of (\mathbf{V}, \mathbf{W}) given $\mathbf{W} = \mathbf{S}(\mathbf{Y})$. Note also that in the latter case one would get directly the marginal likelihood of \mathbf{V} , as \mathbf{V} is independent of \mathbf{W} .

Example 34. Let Y_{ij} , $i \in \{1, \dots, N\}$, $j \in \{1, 2\}$ be independent random variables such that $Y_{i1} \sim \text{Exp}(\psi_i)$ and $Y_{i2} \sim \text{Exp}(\tau \psi_i)$ where $\tau > 0$ and ψ_i are unknown parameters. Show that the distribution of $V_i = \frac{Y_{i2}}{Y_{i1}}$ depends only on parameter τ (and not on ψ_i). Derive the marginal likelihood of τ that is based on $\mathbf{V} = (V_1, \dots, V_N)^\top$.

Example 35. Let Y_{ij} , $i \in \{1, \dots, I\}$, $j \in \{0, 1\}$ be independent, $Y_{ij} \sim \text{Bi}(n_{ij}, p_{ij})$, where

$$\log\left(\frac{p_{ij}}{1-p_{ij}}\right) = \psi_i + \tau \mathbb{I}\{j = 1\}.$$

Suppose we are interested in testing the null hypothesis $H_0 : \tau = 0$ against the alternative $H_1 : \tau \neq 0$.

Note that the standard tests based on the maximum likelihood as described in Chapter 2.6 require that I is fixed and all the sample sizes n_{ij} tend to infinity. This implies that using conditional likelihood is reasonable in situations when (some) n_{ij} are small.

The Rao score test based on the conditional likelihood in this situation coincides with Cochran-Mantel-Haenszel test and its test statistic is given by

$$R_n^{(C)} = \frac{\left(\sum_{i=1}^I Y_{i1} - \mathbf{E}_{H_0}[Y_{i1} | Y_{i+}]\right)^2}{\sum_{i=1}^I \text{var}_{H_0}[Y_{i1} | Y_{i+}]} = \frac{\left(\sum_{i=1}^I Y_{i1} - Y_{i+} \frac{n_{i1}}{n_{i+}}\right)^2}{\sum_{i=1}^I Y_{i+} \frac{n_{i1}n_{i0}}{n_{i+}^2} \frac{n_{i+} - Y_{i+}}{n_{i+} - 1}}, \quad (43)$$

where $Y_{i+} = Y_{i0} + Y_{i1}$ and $n_{i+} = n_{i0} + n_{i1}$. Under the null hypothesis $R_n^{(C)} \xrightarrow[n \rightarrow \infty]{d} \chi_1^2$, where $n = \sum_{i=1}^I \sum_{j=0}^1 n_{ij}$.

Example 36. Consider in Example 35 the special case $I = 1$. Thus the model simplifies to comparing two binomial distributions. Let $Y_0 \sim \text{Bi}(n_0, p_0)$ and $Y_1 \sim \text{Bi}(n_1, p_1)$. Note that the standard approaches of testing the null hypothesis $H_0 : p_0 = p_1$ against the alternative $H_1 : p_0 \neq p_1$ are asymptotic.

The end of
class 10
(22. 3. 2024)

Conditional approach offers an exact inference. Analogously as in Example 35 introduce the parametrization

$$\log\left(\frac{p_j}{1-p_j}\right) = \psi + \tau \mathbb{1}\{j = 1\}, \quad j = 0, 1.$$

Note that in this parametrization τ is the logarithm of odds-ratio.

Put $Y_+ = Y_0 + Y_1$ and $y_+ = y_0 + y_1$. Then

$$P_\tau(Y_1 = k | Y_+ = y_+) = \frac{\binom{n_1}{k} \binom{n_0}{y_+ - k} e^{\tau k}}{\sum_{l \in \mathcal{K}} \binom{n_1}{l} \binom{n_0}{y_+ - l} e^{\tau l}}, \quad k \in \mathcal{K}, \quad (44)$$

where $\mathcal{K} = \{\max\{0, y_+ - n_0\}, \dots, \min\{y_+, n_1\}\}$.

Thus the p-value of the ‘exact’ test of the null hypothesis $H_0 : \tau = \tau_0$ against $H_1 : \tau \neq \tau_0$ would be

$$p(\tau_0) = 2 \min \{ P_{\tau_0}(Y_1 \leq y_1 | Y_+ = y_1 + y_2), P_{\tau_0}(Y_1 \geq y_1 | Y_+ = y_1 + y_2) \}, \quad (45)$$

where y_0 and y_1 are the observed values of Y_0 and Y_1 respectively.

By the inversion of the test one can define the ‘exact’ confidence interval for τ as the set of those values for which we do not reject the null hypothesis, i.e.

$$CI = (\hat{\tau}_L, \hat{\tau}_U) = \{\tau \in \mathbb{R} : p(\tau) > \alpha\}.$$

The confidence interval for odds-ratio calculated by the function `fisher.test()` is now given by $(e^{\hat{\tau}_L}, e^{\hat{\tau}_U})$.

The special case presents testing the null hypothesis $H_0 : \tau = 0$ against $H_1 : \tau \neq 0$. Then (44) simplifies to

$$P_0(Y_1 = k | Y_+ = y_+) = \frac{\binom{n_1}{k} \binom{n_0}{y_+ - k}}{\sum_{l \in \mathcal{K}} \binom{n_1}{l} \binom{n_0}{y_+ - l}} = \frac{\binom{n_1}{k} \binom{n_0}{y_+ - k}}{\binom{n_1 + n_0}{y_+}}, \quad k \in \mathcal{K}.$$

This corresponds to *Fisher’s exact test* sometimes known also as *Fisher’s factorial test*. Be careful that the p-value of the test as implemented in `fisher.test()` is not calculated by (45) but as

$$\tilde{p} = \sum_{k \in \mathcal{K}_-} P_0(Y_1 = k | Y_+ = y_+),$$

where

$$\mathcal{K}_- = \{k \in \mathcal{K} : P_0(Y_1 = k | Y_+ = y_+) \leq P_0(Y_1 = y_1 | Y_+ = y_+)\},$$

which sometimes slightly differs from $p(0)$ as defined in (45).

Note that Fisher’s exact test presents an alternative to the χ^2 -square test of independence in the 2×2 contingency table

	Sample 1	Sample 2
Success	y_0	y_1
Failure	$n_0 - y_0$	$n_1 - y_1$

,

which is an asymptotic test.

Example 37. Consider in Example 35 the special case $n_{i0} = n_{i1} = 1$ for each $i \in \{1, \dots, I\}$.

Introduce

$$N_{jk} = \sum_{i=1}^I \mathbb{1}\{Y_{i0} = j, Y_{i1} = k\}, \quad j, k \in \{0, 1\}.$$

Then the test statistic (43) simplifies to

$$R_n^{(C)} = \frac{(N_{01} - N_{10})^2}{N_{01} + N_{10}},$$

which is known as McNemar's test.

Example 38. Let X_1, \dots, X_{n_1} and Y_1, \dots, Y_{n_2} be two independent random samples from the Poisson distributions. Let λ_X be the true value of the parameter for the first sample and λ_Y for the second sample. Note that $\mathbf{S} = (S_1, S_2)^\top = (\sum_{i=1}^{n_1} X_i, \sum_{i=1}^{n_2} Y_i)^\top$ is a sufficient statistic for the parameter $\boldsymbol{\theta} = (\lambda_X, \lambda_Y)^\top$. Derive the conditional distribution of S_1 given $S_1 + S_2$. Use this result to find an exact test of

$$H_0 : \lambda_X = \lambda_Y, \quad H_1 : \lambda_X \neq \lambda_Y.$$

Further derive an 'exact' confidence interval for the ratio $\frac{\lambda_X}{\lambda_Y}$.

Literature: Pawitan [2001] Chapters 10.1–10.5.

The end of
class 11
(26. 3. 2024)

3 M - and Z -estimators

M - and Z -estimators present a very general class of estimators that include most of the commonly used estimators.*

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from a distribution F from the model \mathcal{F} and one is interested in estimating some quantity (p -dimensional parameter) of this distribution, say $\boldsymbol{\theta}(F)$.

* One should be careful as the terminology may vary. But (among others) minimum contrast estimators, pseudo-likelihood estimators, quasi-likelihood estimators and estimating equations can be usually viewed as either M -estimators or Z -estimators.

***M*-estimator**

Let ρ be a function defined on $S_{\mathbf{X}} \times \Theta$, where $S_{\mathbf{X}}$ is the support of F . Further denote Θ the parameter space, i.e. $\Theta = \{\boldsymbol{\theta}(F), F \in \mathcal{F}\}$. The *M*-estimator* is defined as

$$\hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{X}_i; \boldsymbol{\theta}).$$

Examples of M-estimators

Note that in *parametric models* $\{f(\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ the maximum likelihood (ML) estimator

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \log f(\mathbf{X}_i; \boldsymbol{\theta})$$

can be viewed as the *M*-estimator with

$$\rho_{ML}(\mathbf{x}; \boldsymbol{\theta}) = -\log f(\mathbf{x}; \boldsymbol{\theta}).$$

In *regression problems* one observes $\mathbf{Z}_1 = \begin{pmatrix} Y_1 \\ \mathbf{X}_1 \end{pmatrix}, \dots, \mathbf{Z}_n = \begin{pmatrix} Y_n \\ \mathbf{X}_n \end{pmatrix}$. Then the least squares (LS) estimator of regression parameters

$$\hat{\boldsymbol{\beta}}_n^{(LS)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \mathbf{b})^2$$

can be viewed as the *M*-estimator with

$$\rho_{LS}(\mathbf{z}; \boldsymbol{\beta}) = \rho_{LS}(\mathbf{x}, y; \boldsymbol{\beta}) = (y - \mathbf{x}^T \boldsymbol{\beta})^2.$$

Similarly the least absolute deviation (LAD) estimator

$$\hat{\boldsymbol{\beta}}_n^{(LAD)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n |Y_i - \mathbf{X}_i^T \mathbf{b}|$$

can be viewed as the *M*-estimator with

$$\rho_{LAD}(\mathbf{z}; \boldsymbol{\beta}) = \rho(\mathbf{x}, y; \boldsymbol{\beta}) = |y - \mathbf{x}^T \boldsymbol{\beta}|.$$

***Z*-estimator**

Often the maximizing value in the definition of *M*-estimator is sought by setting a derivative (or the set of partial derivatives if $\boldsymbol{\theta}$ is multidimensional) equal to zero. Thus we search for $\hat{\boldsymbol{\theta}}_n$ as the point that solves the set of estimating equations

$$\frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \hat{\boldsymbol{\theta}}_n) = \mathbf{0}_p, \quad \text{where} \quad \boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial \rho(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}. \quad (46)$$

* *M*-odhad

Note that

$$\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = (\psi_1(\mathbf{x}; \boldsymbol{\theta}), \dots, \psi_p(\mathbf{x}; \boldsymbol{\theta}))^\top = \left(\frac{\partial \rho(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_1}, \dots, \frac{\partial \rho(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_p} \right)^\top.$$

Generally let $\boldsymbol{\psi}$ be a p -dimensional vector function (not necessarily a derivative of the function ρ) defined on $S_{\mathbf{X}} \times \Theta$. Then we define the Z -estimator* as the solution of the system of equations (46).

Note that the maximum likelihood (ML) and the least squares (LS) estimators can be also viewed as Z -estimators with

$$\boldsymbol{\psi}_{ML}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad \boldsymbol{\psi}_{LS}(\mathbf{x}, y; \boldsymbol{\beta}) = (y - \mathbf{x}^\top \boldsymbol{\beta}) \mathbf{x}.$$

Literature: van der Vaart [2000] – Chapter 5.1.

3.1 Identifiability of parameters[†] via M - and/or Z -estimators

When using M - or Z -estimators one should check the potential of these estimators to identify the parameters of interest. Note that by the law of large numbers

$$\frac{1}{n} \sum_{i=1}^n \rho(\mathbf{X}_i; \boldsymbol{\theta}) = \mathbf{E} \rho(\mathbf{X}_1; \boldsymbol{\theta}) + o_P(1), \quad \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}) = \mathbf{E} \boldsymbol{\psi}(\mathbf{X}_1; \boldsymbol{\theta}) + o_P(1).$$

Thus the M -estimator $\hat{\boldsymbol{\theta}}_n$ identifies (at the population level) the quantity

$$\boldsymbol{\theta}_X = \arg \min_{\boldsymbol{\theta} \in \Theta} \mathbf{E} \rho(\mathbf{X}_1; \boldsymbol{\theta})$$

and analogously Z -estimator identifies $\boldsymbol{\theta}_X$ such that

$$\mathbf{E} \boldsymbol{\psi}(\mathbf{X}_1; \boldsymbol{\theta}_X) = \mathbf{0}_p.$$

Example 39. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. observations from a distribution with a density $f(\mathbf{x})$ (with respect to a σ -finite measure μ). By assuming that f belongs to a parametric family of densities $\mathcal{F} = \{f(\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ we are estimating (identifying) $\boldsymbol{\theta}_X$ such that

$$\boldsymbol{\theta}_X = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbf{E} \log f(\mathbf{X}_1; \boldsymbol{\theta})$$

Provided that the true density $f(\mathbf{x})$ has the support $S_{\mathbf{X}}$ that is the same as the support of $f(\mathbf{x}; \boldsymbol{\theta})$ for each $\boldsymbol{\theta} \in \Theta$, this can be further rewritten as

$$\boldsymbol{\theta}_X = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbf{E} \log \left[\frac{f(\mathbf{X}_1; \boldsymbol{\theta})}{f(\mathbf{X}_1)} \right].$$

* Z -odhad † Identifikovatelnost parametru.

Now by Jensen's inequality

$$\mathbb{E} \log \left[\frac{f(\mathbf{X}_1; \boldsymbol{\theta})}{f(\mathbf{X}_1)} \right] \leq \log \left\{ \mathbb{E} \left[\frac{f(\mathbf{X}_1; \boldsymbol{\theta})}{f(\mathbf{X}_1)} \right] \right\} = \log \left\{ \int_{S_{\mathbf{X}}} \frac{f(\mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{x})} f(\mathbf{x}) \, d\mu(\mathbf{x}) \right\} = \log\{1\} = 0.$$

Suppose that our (parametric) assumption **is right** and there exists $\boldsymbol{\theta}_0 \in \Theta$ such that $f(\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}_0)$. Then $\mathbb{E} \log \left[\frac{f(\mathbf{X}_1; \boldsymbol{\theta})}{f(\mathbf{X}_1; \boldsymbol{\theta}_0)} \right]$ is maximised for $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ and thus $\boldsymbol{\theta}_X = \boldsymbol{\theta}_0$ (i.e. the maximum likelihood method identifies the true value of the parameter).

Suppose that our (parametric) assumption **is not right** and that $f \notin \mathcal{F}$. Then

$$\begin{aligned} \boldsymbol{\theta}_X &= \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E} \log \left[\frac{f(\mathbf{X}_1; \boldsymbol{\theta})}{f(\mathbf{X}_1)} \right] = \arg \max_{\boldsymbol{\theta} \in \Theta} \int_{S_{\mathbf{X}}} \log \left[\frac{f(\mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{x})} \right] f(\mathbf{x}) \, d\mu(\mathbf{x}) \\ &= \arg \min_{\boldsymbol{\theta} \in \Theta} \int_{S_{\mathbf{X}}} \log \left[\frac{f(\mathbf{x})}{f(\mathbf{x}; \boldsymbol{\theta})} \right] f(\mathbf{x}) \, d\mu(\mathbf{x}). \end{aligned}$$

The integral $\int_{S_{\mathbf{X}}} \log \left[\frac{f(\mathbf{x})}{f(\mathbf{x}; \boldsymbol{\theta})} \right] f(\mathbf{x}) \, d\mu(\mathbf{x})$ is called the *Kullback–Leibler divergence* from $f(\mathbf{x}; \boldsymbol{\theta})$ to $f(\mathbf{x})$ (it measures how $f(\mathbf{x}; \boldsymbol{\theta})$ diverges from $f(\mathbf{x})$). Thus $\boldsymbol{\theta}_X$ is the point of parameter space Θ for which the Kullback–Leibler divergence from \mathcal{F} to f is minimised.

3.2 Asymptotic distribution of Z -estimators

Analogously as for the maximum likelihood estimator the basic asymptotic results will be formulated for Z -estimators. In order to do that put $\mathbf{Z}(\boldsymbol{\theta}) = \mathbb{E} \boldsymbol{\psi}(\mathbf{X}_1; \boldsymbol{\theta})$ and $\mathbb{D}_{\boldsymbol{\psi}}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial \boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top}$ (the Jacobi matrix of $\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$).

To state the theorem about asymptotic normality we will need the following regularity assumptions. These assumptions are analogous to assumptions **[R0]**–**[R6]** for the maximum likelihood estimators.

[Z0] *Identifiability.* $\boldsymbol{\theta}_X$ satisfies $\mathbf{Z}(\boldsymbol{\theta}_X) = \mathbf{0}_p$.

[Z1] The number of parameters p in the model is *constant*.

[Z2] (The true value of the parameter) $\boldsymbol{\theta}_X$ is an *interior point* of the parameter space Θ .

[Z3] For μ -almost all \mathbf{x} each component of the function $\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta})$ is *differentiable* with respect to $\boldsymbol{\theta}$.

[Z4] There exists $\alpha > 0$, an open neighbourhood U of $\boldsymbol{\theta}_X$ and a function $M(\mathbf{x})$ so that for each $j, k \in \{1, \dots, p\}$ and for each $\boldsymbol{\theta} \in U$

$$\left| \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_k} - \frac{\partial \psi_j(\mathbf{x}; \boldsymbol{\theta}_X)}{\partial \theta_k} \right| \leq M(\mathbf{x}) \|\boldsymbol{\theta} - \boldsymbol{\theta}_X\|^\alpha$$

for μ -almost all \mathbf{x} and $\mathbb{E} M(\mathbf{X}_1) < \infty$.

[Z5] The matrix

$$\mathbb{F}(\boldsymbol{\theta}_X) = \mathbb{E} \mathbb{D}_\psi(\mathbf{X}_1; \boldsymbol{\theta}_X) \quad (47)$$

is finite and *regular*.

[Z6] The *variance matrix*

$$\mathbb{\Sigma}(\boldsymbol{\theta}_X) = \text{var}(\boldsymbol{\psi}(\mathbf{X}_1; \boldsymbol{\theta}_X)) = \mathbb{E} \left[\boldsymbol{\psi}(\mathbf{X}_1; \boldsymbol{\theta}_X) \boldsymbol{\psi}^\top(\mathbf{X}_1; \boldsymbol{\theta}_X) \right] \quad (48)$$

is finite.

Introduce

$$\mathbb{F}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbb{D}_\psi(\mathbf{X}_i; \boldsymbol{\theta}).$$

The following technical lemma says that if $\boldsymbol{\theta}$ is ‘close’ to $\boldsymbol{\theta}_X$, then $\mathbb{F}_n(\boldsymbol{\theta})$ is close to $\mathbb{F}(\boldsymbol{\theta}_X)$. This result will be useful for the proof of the consistency and asymptotic normality of Z -estimators. Note that it is an analogy of Lemma 1.

Lemma 4. *Suppose that assumptions [Z1]-[Z5] are satisfied. Let $\{\varepsilon_n\}$ be a sequence of positive numbers going to zero. Then*

$$\max_{j,k \in \{1, \dots, p\}} \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (\mathbb{F}_n(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}_X))_{jk} \right| = o_P(1),$$

where

$$U_{\varepsilon_n} = \{ \boldsymbol{\theta} \in \Theta : \|\boldsymbol{\theta} - \boldsymbol{\theta}_X\| \leq \varepsilon_n \}$$

and $(\mathbb{F}_n(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}_X))_{jk}$ stands for the (j, k) -element of the difference of the matrices $\mathbb{F}_n(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}_X)$.

Proof. Using assumption [Z4] and the law of large numbers one can bound

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (\mathbb{F}_n(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}_X))_{jk} \right| &\leq \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (\mathbb{F}_n(\boldsymbol{\theta}) - \mathbb{F}_n(\boldsymbol{\theta}_X))_{jk} \right| + \left| (\mathbb{F}_n(\boldsymbol{\theta}_X) - \mathbb{F}(\boldsymbol{\theta}_X))_{jk} \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n M(\mathbf{X}_i) \varepsilon_n^\alpha + o_P(1) = O_P(1) o(1) + o_P(1) = o_P(1), \end{aligned}$$

which implies the statement of the lemma. \square

Suppose now that $\widehat{\mathbf{t}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$. Note that the above lemma (together with the reasoning of Corollary 1) one gets that for each $j, k \in \{1, \dots, p\}$:

$$\left| (\mathbb{F}_n(\widehat{\mathbf{t}}_n) - \mathbb{F}(\boldsymbol{\theta}_X))_{jk} \right| = o_P(1). \quad (49)$$

Theorem 9. *Suppose that assumptions [Z0]-[Z6] are satisfied.*

(i) Then with probability going to one there exists a consistent solution $\widehat{\boldsymbol{\theta}}_n$ to the estimating equations (46).

(ii) Further, if $\widehat{\boldsymbol{\theta}}_n$ is a consistent root of the estimating equations (46), then

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = -\mathbb{F}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X) + o_P(1), \quad (50)$$

which further implies that

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{F}^{-1}(\boldsymbol{\theta}_X) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{F}^{-1}(\boldsymbol{\theta}_X)]^\top), \quad (51)$$

where the matrices $\mathbb{F}(\boldsymbol{\theta}_X)$ and $\boldsymbol{\Sigma}(\boldsymbol{\theta}_X)$ are defined in (47) and (48) respectively.

Proof. Consistency: Introduce the vector function

$$h_n(\boldsymbol{\theta}) = \boldsymbol{\theta} - [\mathbb{F}(\boldsymbol{\theta}_X)]^{-1} \mathbf{Z}_n(\boldsymbol{\theta}),$$

where

$$\mathbf{Z}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}).$$

In what follows we will show that with probability going to one (as $n \rightarrow \infty$) the mapping h_n is a contraction on $U_{\varepsilon_n} = \{\boldsymbol{\theta} \in \Theta : \|\boldsymbol{\theta} - \boldsymbol{\theta}_X\| \leq \varepsilon_n\}$, where $\{\varepsilon_n\}$ is a sequence of positive numbers going to zero such that $\varepsilon_n \sqrt{n} \xrightarrow[n \rightarrow \infty]{} \infty$. Having proved that then by the Banach fixed point theorem (Theorem A14) there exists a unique fixed point $\widehat{\boldsymbol{\theta}}_n \in U_{\varepsilon_n}$ such that $h_n(\widehat{\boldsymbol{\theta}}_n) = \widehat{\boldsymbol{\theta}}_n$ and thus also $\mathbf{Z}_n(\widehat{\boldsymbol{\theta}}_n) = \mathbf{0}_p$. This implies the existence of a consistent root of the estimating equations (46).

Showing that h_n is a contraction on U_{ε_n} . Let $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in U_{\varepsilon_n}$ then

$$\begin{aligned} \|h_n(\boldsymbol{\theta}_1) - h_n(\boldsymbol{\theta}_2)\| &= \|(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2) - [\mathbb{F}(\boldsymbol{\theta}_X)]^{-1}(\mathbf{Z}_n(\boldsymbol{\theta}_1) - \mathbf{Z}_n(\boldsymbol{\theta}_2))\| \\ &= \|(\mathbb{I}_p - [\mathbb{F}(\boldsymbol{\theta}_X)]^{-1} \mathbb{F}_n^*)(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)\|, \end{aligned} \quad (52)$$

where \mathbb{F}_n^* is $(p \times p)$ -matrix whose j -th row is the j -th row of the matrix

$$\mathbb{F}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbb{D}_{\boldsymbol{\psi}}(\mathbf{X}_i; \boldsymbol{\theta})$$

evaluated at some $\boldsymbol{\theta}_n^{j*}$ that is between $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$. Note that $\boldsymbol{\theta}_n^{j*} \in U_{\varepsilon_n}$. Now by Lemma 4 and assumption [Z5]

$$a_n = \max_{j,k \in \{1, \dots, p\}} \sup_{\boldsymbol{\theta} \in U_{\varepsilon_n}} \left| (\mathbb{I}_p - [\mathbb{F}(\boldsymbol{\theta}_X)]^{-1} \mathbb{F}_n(\boldsymbol{\theta}))_{jk} \right| = o_P(1). \quad (53)$$

So with the help of (52) and (53) it holds that uniformly in $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in U_{\varepsilon_n}$

$$\|h_n(\boldsymbol{\theta}_1) - h_n(\boldsymbol{\theta}_2)\| \leq o_P(1) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|, \quad (54)$$

which implies that there exists $q \in (0, 1)$ such that

$$\mathbb{P}\left(\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in U_{\varepsilon_n} \ \|h_n(\boldsymbol{\theta}_1) - h_n(\boldsymbol{\theta}_2)\| \leq q \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|\right) \xrightarrow{n \rightarrow \infty} 1.$$

Thus to show that h_n is a contraction on U_{ε_n} it remains to prove that (with probability going to one) $h_n : U_{\varepsilon_n} \rightarrow U_{\varepsilon_n}$. Note that for each $\boldsymbol{\theta} \in U_{\varepsilon_n}$ the inequality (54) implies

$$h_n(\boldsymbol{\theta}) - h_n(\boldsymbol{\theta}_X) = o_P(1) \varepsilon_n, \quad (55)$$

where the $o_P(1)$ term does not depend on $\boldsymbol{\theta}$. Further

$$h_n(\boldsymbol{\theta}_X) = \boldsymbol{\theta}_X - [\mathbb{F}(\boldsymbol{\theta}_X)]^{-1} \mathbf{Z}_n(\boldsymbol{\theta}_X) = \boldsymbol{\theta}_X + O_P\left(\frac{1}{\sqrt{n}}\right), \quad (56)$$

where we have used that by the central limit theorem $\mathbf{Z}_n(\boldsymbol{\theta}_X) = O_P\left(\frac{1}{\sqrt{n}}\right)$. Now combining (55) and (56) yields that

$$h_n(\boldsymbol{\theta}) = o_P(1) \varepsilon_n + h_n(\boldsymbol{\theta}_X) = o_P(1) \varepsilon_n + \boldsymbol{\theta}_X + O_P\left(\frac{1}{\sqrt{n}}\right),$$

which further together with the assumption $\varepsilon_n \sqrt{n} \xrightarrow{n \rightarrow \infty} \infty$ gives

$$h_n(\boldsymbol{\theta}) - \boldsymbol{\theta}_X = \varepsilon_n \left(o_P(1) + O_P\left(\frac{1}{\varepsilon_n \sqrt{n}}\right) \right) = \varepsilon_n o_P(1).$$

This finally implies that $\mathbb{P}(\forall \boldsymbol{\theta} \in U_{\varepsilon_n} : h_n(\boldsymbol{\theta}) \in U_{\varepsilon_n}) \xrightarrow{n \rightarrow \infty} 1$, which was to be proved.

Asymptotic normality: This is proved analogously as in Theorem 5. Let $\widehat{\boldsymbol{\theta}}_n$ be a consistent root of the estimating equations. Then by the mean value theorem applied to each component of $\mathbf{Z}_n(\widehat{\boldsymbol{\theta}}_n)$ one gets

$$\mathbf{0}_p = \mathbf{Z}_n(\widehat{\boldsymbol{\theta}}_n) = \mathbf{Z}_n(\boldsymbol{\theta}_X) + \mathbb{F}_n^*(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X),$$

where similarly as in the proof of consistency \mathbb{F}_n^* is $(p \times p)$ -matrix whose j -th row is the j -th row of the matrix $\mathbb{F}_n(\boldsymbol{\theta})$ evaluated at some $\boldsymbol{\theta}_n^{j*}$ that is between $\widehat{\boldsymbol{\theta}}_n$ and $\boldsymbol{\theta}_X$. Thus $\boldsymbol{\theta}_n^{j*} \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\theta}_X$ as $\widehat{\boldsymbol{\theta}}_n$ is a consistent estimator of $\boldsymbol{\theta}_X$. So one can use (49) to conclude that $\mathbb{F}_n^* \xrightarrow[n \rightarrow \infty]{P} \mathbb{F}(\boldsymbol{\theta}_X)$. Now with the help of CS (Theorem 2) one can write

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = -[\mathbb{F}_n^*]^{-1} \sqrt{n} \mathbf{Z}_n(\boldsymbol{\theta}_X) = -\mathbb{F}^{-1}(\boldsymbol{\theta}_X) \frac{1}{\sqrt{n}} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X) + o_P(1),$$

which with the help of the central limit theorem (for i.i.d. random vectors) and CS (Theorem 2) implies the second statement of the theorem. \square

Remark 15. If there exists a real function $\rho(\mathbf{x}; \boldsymbol{\theta})$ such that $\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\partial \rho(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$, then the matrix $\mathbb{F}(\boldsymbol{\theta}_X)$ is symmetric and one can simply write $\mathbb{F}(\boldsymbol{\theta}_X)^{-1}$ instead of $[\mathbb{F}(\boldsymbol{\theta}_X)^{-1}]^\top$ in (51).

Asymptotic variance estimations

Note that by Theorem 9 one has

$$\widehat{\boldsymbol{\theta}}_n \stackrel{\text{as}}{\approx} \mathbf{N}_p(\boldsymbol{\theta}_X, \frac{1}{n} \mathbb{F}^{-1}(\boldsymbol{\theta}_X) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{F}^{-1}(\boldsymbol{\theta}_X)]^\top).$$

Thus the most straightforward estimate of the asymptotic variance of $\widehat{\boldsymbol{\theta}}_n$ is the ‘sandwich estimator’ given by

$$\widehat{\text{avar}}(\widehat{\boldsymbol{\theta}}_n) = \frac{1}{n} \widehat{\mathbb{F}}_n^{-1} \widehat{\boldsymbol{\Sigma}}_n [\widehat{\mathbb{F}}_n^{-1}]^\top, \quad (57)$$

where

$$\widehat{\mathbb{F}}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{D}\boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) \quad \text{and} \quad \widehat{\boldsymbol{\Sigma}}_n = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) \boldsymbol{\psi}^\top(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n).$$

Note that Lemma 4 together with the consistency of $\widehat{\boldsymbol{\theta}}_n$ implies that

$$\widehat{\mathbb{F}}_n \xrightarrow[n \rightarrow \infty]{P} \mathbb{F}(\boldsymbol{\theta}_X).$$

It is more tedious to give some general assumptions so that it also holds

$$\widehat{\boldsymbol{\Sigma}}_n \xrightarrow[n \rightarrow \infty]{P} \boldsymbol{\Sigma}(\boldsymbol{\theta}_X).$$

To derive such assumptions rewrite

$$\begin{aligned} \widehat{\boldsymbol{\Sigma}}_n &= \frac{1}{n} \sum_{i=1}^n [\boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) - \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X)] [\boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) - \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X)]^\top \\ &\quad + \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X) [\boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) - \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X)]^\top \\ &\quad + \frac{1}{n} \sum_{i=1}^n [\boldsymbol{\psi}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) - \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X)] \boldsymbol{\psi}^\top(\mathbf{X}_i; \boldsymbol{\theta}_X) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}(\mathbf{X}_i; \boldsymbol{\theta}_X) \boldsymbol{\psi}^\top(\mathbf{X}_i; \boldsymbol{\theta}_X). \end{aligned} \quad (58)$$

Now by the law of large numbers the last summand in (58) converges in probability to $\boldsymbol{\Sigma}(\boldsymbol{\theta}_X)$, thus it is sufficient to show that the remaining terms are of order $o_P(1)$. With the help of assumption [Z4] this can be done for instance by assuming that for each $j, k \in \{1, \dots, p\}$

$$\mathbb{E} M_{jk}^2(\mathbf{X}_1) < \infty \quad \text{and} \quad \mathbb{E} \left| \frac{\partial \psi_j(\mathbf{X}_1; \boldsymbol{\theta}_X)}{\partial \theta_k} \right|^2 < \infty.$$

Confidence sets and confidence intervals

Suppose that $\widehat{\mathbb{V}}_n$ is a consistent estimator of $\mathbb{V} = \mathbb{F}^{-1}(\boldsymbol{\theta}_X) \boldsymbol{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{F}^{-1}(\boldsymbol{\theta}_X)]^\top$.

Then by the Cramér-Slutsky theorem the confidence set (ellipsoid) for the parameter $\boldsymbol{\theta}_X$ is given by

$$\left\{ \boldsymbol{\theta} \in \Theta : n (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta})^\top \widehat{\mathbb{V}}_n^{-1} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \leq \chi_p^2(1 - \alpha) \right\}.$$

The ‘Wald-type’ (asymptotic) confidence interval for θ_{Xj} (the j -th coordinate of $\boldsymbol{\theta}_X$) is given by

$$\left[\widehat{\theta}_{nj} - \frac{u_{1-\alpha/2} \sqrt{\widehat{v}_{n,jj}}}{\sqrt{n}}, \widehat{\theta}_{nj} + \frac{u_{1-\alpha/2} \sqrt{\widehat{v}_{n,jj}}}{\sqrt{n}} \right],$$

where $\widehat{\theta}_{nj}$ is the j -th coordinate of $\widehat{\boldsymbol{\theta}}_n$ and $\widehat{v}_{n,jj}$ is the j -th diagonal element of the matrix $\widehat{\mathbb{V}}_n$.

Literature: Sen et al. [2010] Chapter 8.2.

3.3 Likelihood under model misspecification

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample with a density f (with respect to a σ -finite measure μ). Then the maximum likelihood estimator can be viewed as the M -estimator with $\rho(\mathbf{x}; \boldsymbol{\theta}) = -\log f(\mathbf{x}; \boldsymbol{\theta})$ or Z -estimator with $\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = -\frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$. From Example 39 we know that when assuming $f \in \mathcal{F} = \{f(\mathbf{x}; \boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta\}$, the method of the maximum likelihood identifies the parameter

$$\boldsymbol{\theta}_X = \arg \min_{\boldsymbol{\theta} \in \Theta} \int_{S_{\mathbf{X}}} \log \left[\frac{f(\mathbf{x})}{f(\mathbf{x}; \boldsymbol{\theta})} \right] f(\mathbf{x}) d\mu(\mathbf{x}).$$

Further by Theorem 9 we also know that (with probability going to one there exists a consistent solution $\widehat{\boldsymbol{\theta}}_n$ of (46) which satisfies)

$$\sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{F}^{-1}(\boldsymbol{\theta}_X) \mathbb{\Sigma}(\boldsymbol{\theta}_X) [\mathbb{F}^{-1}(\boldsymbol{\theta}_X)]^\top).$$

Suppose that our parametric assumption is right and $f \in \mathcal{F}$, i.e. there exists $\boldsymbol{\theta}_0 \in \Theta$ such that $f(\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}_0)$. Then the identified parameter is equal to $\boldsymbol{\theta}_0$, i.e. $\boldsymbol{\theta}_X = \boldsymbol{\theta}_0$. Further it is easy to see that $\mathbb{F}(\boldsymbol{\theta}_X) = I(\boldsymbol{\theta}_X) = \mathbb{\Sigma}(\boldsymbol{\theta}_X)$, where $I(\boldsymbol{\theta}_X)$ is the Fisher information matrix. Thus

$$\mathbb{F}^{-1}(\boldsymbol{\theta}_X) \mathbb{\Sigma}(\boldsymbol{\theta}_X) \mathbb{F}^{-1}(\boldsymbol{\theta}_X) = I^{-1}(\boldsymbol{\theta}_X).$$

So one can view Theorem 5 as a special case of Theorem 9. Further, when doing the inference about $\boldsymbol{\theta}_X$ it is sufficient to estimate the Fisher information matrix.

Often in practice we are not completely sure that $f \in \mathcal{F}$. If we are not sure about the parametric assumption then it is safer to view the estimator $\widehat{\boldsymbol{\theta}}_n$ as an Z -estimator with $\boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}) = -\frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$. The asymptotic variance of $\widehat{\boldsymbol{\theta}}_n$ can now be estimated with the help of ‘sandwich estimator’ (57) where

$$\begin{aligned} \widehat{\mathbb{\Sigma}}_n &= \frac{1}{n} \sum_{i=1}^n \mathbf{U}(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n) \mathbf{U}^\top(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n), \text{ where } \mathbf{U}(\mathbf{x}; \boldsymbol{\theta}) = -\frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \\ \widehat{\mathbb{F}}_n &= \frac{1}{n} \sum_{i=1}^n I(\mathbf{X}_i; \widehat{\boldsymbol{\theta}}_n), \text{ where } I(\mathbf{x}; \boldsymbol{\theta}) = -\frac{\partial^2 \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}. \end{aligned}$$

This type of variance estimator is calculated for GLM models by the function `sandwich` (from the package with the same name).

Example 40. *Misspecified normal linear model.* Let $(\begin{smallmatrix} Y_1 \\ \mathbf{X}_1 \end{smallmatrix}), \dots, (\begin{smallmatrix} Y_n \\ \mathbf{X}_n \end{smallmatrix})$ be independent and identically distributed random vectors, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$. Note that if one assumes that $\mathcal{L}(Y_i|\mathbf{X}_i) \sim \mathbf{N}(\mathbf{X}_i^\top \boldsymbol{\beta}, \sigma^2)$ for some $\boldsymbol{\beta} \in \mathbb{R}^p$, then the maximum likelihood estimation of $\boldsymbol{\beta}$ corresponds to the method of the least squares given by $\rho_{LS}(\mathbf{x}, y; \boldsymbol{\beta}) = (y - \mathbf{x}^\top \boldsymbol{\beta})^2$.

Show that without the assumption $\mathcal{L}(Y_i|\mathbf{X}_i) \sim \mathbf{N}(\mathbf{X}_i^\top \boldsymbol{\beta}, \sigma^2)$ the method of the least squares identifies the parameter

$$\boldsymbol{\beta}_X = [\mathbf{E} \mathbf{X}_1 \mathbf{X}_1^\top]^{-1} \mathbf{E} Y_1 \mathbf{X}_1$$

and it holds that $\sqrt{n} (\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{V})$, where

$$\mathbb{V} = [\mathbf{E} \mathbf{X}_1 \mathbf{X}_1^\top]^{-1} [\mathbf{E} \sigma^2(\mathbf{X}_1) \mathbf{X}_1 \mathbf{X}_1^\top] [\mathbf{E} \mathbf{X}_1 \mathbf{X}_1^\top]^{-1},$$

with $\sigma^2(\mathbf{X}_1) = \mathbf{E} [(Y_1 - \mathbf{X}_1^\top \boldsymbol{\beta}_X)^2 | \mathbf{X}_1]$.

Note that provided $\mathbf{E} [Y_1 | \mathbf{X}_1] = \mathbf{X}_1^\top \boldsymbol{\beta}_0$ for some $\boldsymbol{\beta}_0 \in \mathbb{R}^p$, then $\boldsymbol{\beta}_X = \boldsymbol{\beta}_0$ and $\sigma^2(\mathbf{X}_1) = \text{var}(Y_1 | \mathbf{X}_1)$.

Example 41. *Misspecified Poisson regression.* Let $(\begin{smallmatrix} Y_1 \\ \mathbf{X}_1 \end{smallmatrix}), \dots, (\begin{smallmatrix} Y_n \\ \mathbf{X}_n \end{smallmatrix})$ be independent and identically distributed random vectors, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$. Assume that the conditional distribution of Y_i given \mathbf{X}_i is Poisson, i.e. $\mathcal{L}(Y_i|\mathbf{X}_i) \sim \text{Po}(\lambda(\mathbf{X}_i))$, where $\lambda(\mathbf{x}) = e^{\mathbf{x}^\top \boldsymbol{\beta}}$ and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$. The score statistic for the maximum likelihood estimation is given by

$$\mathbf{U}_n(\boldsymbol{\beta}) = \sum_{i=1}^n \mathbf{X}_i (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}}).$$

Thus one can view the maximum likelihood estimator $\widehat{\boldsymbol{\beta}}_n$ as the Z -estimator with

$$\boldsymbol{\psi}(\mathbf{x}, y; \boldsymbol{\beta}) = \mathbf{x} (y - e^{\mathbf{x}^\top \boldsymbol{\beta}}) \tag{59}$$

and $\boldsymbol{\beta}_X$ solves the system of equations

$$\mathbf{E} \mathbf{X}_1 (Y_1 - e^{\mathbf{X}_1^\top \boldsymbol{\beta}_X}) = \mathbf{0}_p.$$

Suppose now that $\mathcal{L}(Y_i|\mathbf{X}_i) \not\sim \text{Po}(\lambda(\mathbf{X}_i))$, but one can still assume that there exists $\boldsymbol{\beta}_0$ such that $\mathbf{E} [Y_1 | \mathbf{X}_1] = e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}$. Then

$$\mathbf{E} \mathbf{X}_1 (Y_1 - e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}) = \mathbf{E} \left\{ \mathbf{E} [\mathbf{X}_1 (Y_1 - e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}) | \mathbf{X}_1] \right\} = \mathbf{E} [\mathbf{X}_1 (e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0} - e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0})] = \mathbf{0}_p.$$

Thus $\boldsymbol{\beta}_X$ identifies $\boldsymbol{\beta}_0$ which describes the effect of the covariates on the expected mean value.

The above calculation implies that when we are not sure that the conditional distribution $\mathcal{L}(Y_i|\mathbf{X}_i)$ is $\text{Po}(\lambda(\mathbf{X}_i))$, but we are willing to assume that $\mathbf{E} [Y_i | \mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ for some $\boldsymbol{\beta}_0 \in \mathbb{R}^p$,

then we can still use the score function (59) which identifies the parameter β_0 . By Theorem 9 we know that the estimator $\widehat{\beta}_n$ is asymptotically normal with the matrices $\mathbb{F}(\beta_X)$ and $\mathbb{Z}(\beta_X)$ given by

$$\mathbb{Z}(\beta_X) = \mathbb{E} \mathbf{X}_1 \mathbf{X}_1^\top (Y_1 - e^{\mathbf{X}_1^\top \beta_X})^2 \quad \text{and} \quad \mathbb{F}(\beta_X) = \mathbb{E} \mathbf{X}_1 \mathbf{X}_1^\top e^{\mathbf{X}_1^\top \beta_X}.$$

Thus the asymptotic variance of the estimator $\widehat{\beta}_n$ can be estimated by

$$\widehat{\text{avar}}(\widehat{\beta}_n) = \frac{1}{n} \widehat{\mathbb{F}}_n^{-1} \widehat{\mathbb{Z}}_n \widehat{\mathbb{F}}_n^{-1},$$

where

$$\widehat{\mathbb{Z}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top (Y_i - e^{\mathbf{X}_i^\top \widehat{\beta}_n})^2 \quad \text{and} \quad \widehat{\mathbb{F}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top e^{\mathbf{X}_i^\top \widehat{\beta}_n}.$$

Literature: [White \[1980\]](#), [White \[1982\]](#).

3.4 Asymptotic normality of M -estimators defined by convex minimization

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from a distribution F and one is interested in estimating some quantity θ_X (p -dimensional parameter) of this distribution such that this parameter can be identified as

$$\theta_X = \arg \min_{\theta \in \Theta} \mathbb{E} \rho(\mathbf{X}_1; \theta),$$

where for each fixed \mathbf{x} the function $\rho(\mathbf{x}; \theta)$ is convex in θ . Further suppose that the parameter space Θ is a subset of \mathbb{R}^p .

A straightforward estimator of the parameter θ_X is given by

$$\widehat{\theta}_n = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{X}_i; \theta).$$

As we see later we do not need to assume that the function $\rho(\mathbf{x}; \theta)$ is differentiable in θ for all $(\mathbf{x}, \theta) \in S_X \times \Theta$. Nevertheless we need a function that plays the role of the function $\psi(\mathbf{x}; \theta)$ in the definition of Z -estimators (46). Note that the convexity in θ guarantees that for each \mathbf{x} the function $\rho(\mathbf{x}, \theta)$ is differentiable in θ for almost all $\theta \in \Theta$. So suppose that there exists a function $\psi(\mathbf{x}; \theta)$ such that $\psi(\mathbf{x}; \theta) = \frac{\partial \rho(\mathbf{x}; \theta)}{\partial \theta}$ whenever this derivative exists. Moreover suppose that similarly as for Z -estimators it holds that

$$\mathbb{E} \psi(\mathbf{X}_1; \theta_X) = \mathbf{0}_p. \tag{60}$$

For formulating the main result it is useful to introduce the ‘remainder function’

$$R(\mathbf{x}; \mathbf{t}) = \rho(\mathbf{x}; \theta_X + \mathbf{t}) - \rho(\mathbf{x}; \theta_X) - \mathbf{t}^\top \psi(\mathbf{x}; \theta_X) \tag{61}$$

and the asymptotic (expected) objective function

$$M(\theta) = \mathbb{E} \rho(\mathbf{X}_1; \theta). \tag{62}$$

Theorem 10. Assume that the function $\psi(\mathbf{x}; \boldsymbol{\theta})$ satisfies (60) and the functions $R(\mathbf{x}; \mathbf{t})$ and $M(\boldsymbol{\theta})$ are defined by (61) and (62) respectively. Further suppose that

(i) there exists a positive definite matrix $\mathbb{J}(\boldsymbol{\theta}_X)$ such that

$$M(\boldsymbol{\theta}_X + \mathbf{t}) = M(\boldsymbol{\theta}_X) + \frac{1}{2} \mathbf{t}^\top \mathbb{J}(\boldsymbol{\theta}_X) \mathbf{t} + o(\|\mathbf{t}\|^2), \text{ as } \mathbf{t} \rightarrow \mathbf{0}_p;$$

(ii) $\text{var}(R(\mathbf{X}_1; \mathbf{t})) = o(\|\mathbf{t}\|^2)$ as $\mathbf{t} \rightarrow \mathbf{0}_p$;

(iii) there exists a finite variance matrix $\Sigma(\boldsymbol{\theta}_X) = \text{var}(\psi(\mathbf{X}_1; \boldsymbol{\theta}_X))$.

Then

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) = -[\mathbb{J}(\boldsymbol{\theta}_X)]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(\mathbf{X}_i; \boldsymbol{\theta}_X) + o_P(1),$$

which further implies that

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, [\mathbb{J}(\boldsymbol{\theta}_X)]^{-1} \Sigma(\boldsymbol{\theta}_X) [\mathbb{J}(\boldsymbol{\theta}_X)]^{-1}).$$

Proof. See the proof of Theorem 2.1 of Hjort and Pollard [2011]. □

Comparison of Theorems 9 and 10

First of all note that Theorem 10 yields the asymptotic normality of the argument of the minimum of $\frac{1}{n} \sum_{i=1}^n \rho(\mathbf{X}_i; \boldsymbol{\theta})$. On the other hand Theorem 9 guarantees asymptotic normality only for a consistent (i.e. an appropriately chosen) root of the estimating equations (46). But in case that there are more roots to the estimating equations (46) it is generally impossible to decide which of the roots is the consistent one.

Further it is worth noting that Theorem 10 allows for $\rho(\mathbf{x}; \boldsymbol{\theta})$ that are ‘less differentiable’. Note that to calculate the matrix $\Gamma(\boldsymbol{\theta}_X)$ one needs the differentiability of the function ψ . On the other hand the matrix $\mathbb{J}(\boldsymbol{\theta}_X)$ can be computed as the Hessian matrix of the function $M(\boldsymbol{\theta}) = \mathbf{E} \rho(\mathbf{X}_1; \boldsymbol{\theta})$ at the point $\boldsymbol{\theta}_X$. Thus the assumption about the smoothness of ψ (i.e. [Z3] and [Z4]) can be replaced with the assumption that function $M(\boldsymbol{\theta})$ is twice continuously differentiable on a neighbourhood of $\boldsymbol{\theta}_X$. So the lack of smoothness of ψ can be compensated with the assumptions on the distribution of \mathbf{X}_1 so that the function $M(\boldsymbol{\theta})$ is sufficiently smooth in $\boldsymbol{\theta}$. See also the application of Theorem 10 to derive the asymptotic normality of the sample median given below.

Vaguely speaking the assumptions of Theorems 10 are milder than assumptions of Theorem 9. More formally if assumptions [Z3]-[Z6] hold then also assumptions (i) and (iii) are satisfied. Further a closer inspection of the proof of Theorem 2.1 of Hjort and Pollard [2011] shows that the remainder term $r_n(s)$ there can be handled with the help of assumptions [Z3] and [Z4] instead of assumption (ii) of Theorem 10.

On the other hand note that if the function ψ does not meet at least assumption [Z3] then it is not any straightforward method how to estimate the matrix $\mathbb{J}(\theta_X)$ which is needed to estimate the asymptotic variance of $\widehat{\theta}_n$.

Sample median

Let X_1, \dots, X_n be independent identically distributed random variables with density $f(y)$ that is positive and continuous in a neighbourhood of median $F^{-1}(0.5)$.

It is well known (see also Lemma 5 and Remark 17 in Chapter 5) that the sample median \widetilde{m}_n can be written as

$$\widetilde{m}_n = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n |X_i - \theta|.$$

Thus one can view \widetilde{m}_n as an M -estimator with $\rho(x; \theta) = |x - \theta|$. For theoretical reasons it is advantageous to consider

$$\rho(x; \theta) = |x - \theta| - |x|,$$

which does not require that $\mathbb{E} |X_1| < \infty$ in order to define $M(\theta) = \mathbb{E} \rho(X; \theta)$. Note that then (see also Lemma 5 in Chapter 5)

$$F^{-1}(0.5) = \arg \min_{\theta \in \mathbb{R}} \mathbb{E} \rho(X; \theta).$$

Now one can use Theorem 10 to derive the asymptotic distribution of \widetilde{m}_n . Introduce

$$\psi(x; \theta) = -\text{sign}(x - \theta)$$

and note that $\psi(x; \theta) = \partial \rho(x; \theta) / \partial \theta$ for $\theta \neq x$. Further it is easy to check that

$$\mathbb{E} \psi(X_1; F^{-1}(0.5)) = 0.$$

Provided that also the other assumptions of Theorem 10 are satisfied it remains to calculate $\Sigma(\theta_X)$ and $\Gamma(\theta_X)$. As $\theta \in \mathbb{R}$ the matrix $\Sigma(\theta_X)$ reduces to

$$\sigma_\psi^2 = \text{var} (\psi(X_1; F^{-1}(0.5))) = 1.$$

Further as $M(\theta) = -\mathbb{E} \int_0^\theta \text{sign}(X_1 - t) dt$ one can interchange the derivative and the integral to get

$$\frac{\partial M(\theta)}{\partial \theta} = -\mathbb{E} [\text{sign}(X_1 - \theta)] = -\mathbb{P}(X_1 > \theta) + \mathbb{P}(X_1 < \theta) = 2F(\theta) - 1.$$

This further implies that (the matrix) $\mathbb{J}(\theta_X)$ reduces to

$$\gamma = \left. \frac{\partial^2 M(\theta)}{\partial \theta^2} \right|_{\theta=F^{-1}(0.5)} = 2f(F^{-1}(0.5)).$$

Finally one gets

$$\sqrt{n} (\tilde{m}_n - F^{-1}(0.5)) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}\left(0, \frac{1}{4 f^2(F^{-1}(0.5))}\right).$$

Note also to estimate the asymptotic variance of \tilde{m}_n one needs to estimate the quantity $f(F^{-1}(0.5))$ which is far from being straightforward.*

Literature: Hjort and Pollard [2011] Section 2A.

4 M -estimators and Z -estimators in robust statistics†

In statistics the word ‘robust’ has basically two meanings.

- (i) We say that a procedure is robust, if it stays (approximately/asymptotically) valid even when some of the assumptions (under which the procedure is derived) are not satisfied. For instance the standard ANOVA F -statistic is robust against the violation of the normality of the observations provided that the variances of all the observations are the same (and finite).
- (ii) People interested in robust statistics say that a procedure is robust, if it is not ‘too much’ influenced by the outlying observations. In what follows we will concentrate on this meaning of the robustness.

One of the standard measures of robustness is the **breakdown point**. Vaguely speaking the breakdown point of an estimator is the smallest percentage of observations that one has to change so that the estimator produces a nonsense value (e.g. $\pm\infty$ for location or regression estimator; 0 or $+\infty$ when estimating the scale).

Let $\hat{\boldsymbol{\theta}}_n$ be an M - or Z -estimator of a parameter $\boldsymbol{\theta}_X$. Note that thanks to Theorems 9 or 10 (under appropriate assumptions) one has the following representation

$$\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X = \frac{1}{n} \sum_{i=1}^n IF(\mathbf{X}_i) + o_P\left(\frac{1}{\sqrt{n}}\right),$$

where $IF(\mathbf{x}) = -\mathbb{F}^{-1}(\boldsymbol{\theta}_X) \boldsymbol{\psi}(\mathbf{x}; \boldsymbol{\theta}_X)$ is called *the influence function*. Thus if one can ignore the remainder term $o_P\left(\frac{1}{\sqrt{n}}\right)$, then changing \mathbf{X}_i to $\mathbf{X}_i + \boldsymbol{\Delta}$ results that the estimates $\hat{\boldsymbol{\theta}}_n$ changes (approximately) by

$$\frac{1}{n} [IF(\mathbf{X}_i + \boldsymbol{\Delta}) - IF(\mathbf{X}_i)].$$

Thus provided that $IF(\mathbf{x})$ is bounded then also this change is bounded (and of order $O\left(\frac{1}{n}\right)$).

Note that the above reasoning was not completely correct as the term $o_P\left(\frac{1}{\sqrt{n}}\right)$ was ignored. Nevertheless it can be proved that (under some mild assumptions excluding ‘singular’ cases) if

* This can be estimated with the help of kernel smoothing methods or circumvented with the help of bootstrap methods. Both methods are included in the course [NMST545 Mathematical Statistics 4](#). † *Robustní statistika*

the function $\psi(\mathbf{x}; \boldsymbol{\theta})$ is bounded then the breakdown point of the associated $M(Z)$ -estimator is $\frac{1}{2}$.

4.1 Robust estimation of location*

Suppose that we observe a random sample X_1, \dots, X_n from a distribution F and we are interested in characterising the location.

Note that for the sample mean $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ it is sufficient to change only one observation to get an arbitrary value of \bar{X}_n .

On the other hand when considering the sample median $\tilde{m}_n = \hat{F}_n^{-1}(0.5)$ then one needs to change at least half of the observations so that one can for instance change the estimator to $\pm\infty$.

When deciding between a sample mean and a sample median one has to take into consideration that if the distribution F is not symmetric then \bar{X}_n and \tilde{m}_n estimate different quantities. But when one can hope that the distribution F is symmetric, then both \bar{X}_n and \tilde{m}_n estimate the centre of the symmetry and one can be interested which of the estimators is more appropriate. By the maximum likelihood theory we know that \bar{X}_n is efficient if F is normal while \tilde{m}_n is asymptotically efficient if F is doubly exponential (i.e. it has a density $f(x) = \frac{1}{2\sigma} \exp\{-\frac{|x-\theta|}{\sigma}\}$).

In robust statistics it is usually assumed that most of our observations follow normal distributions but there are some outlying values. This can be formalised by assuming that the distribution function F of each of the observations satisfies

$$F(x) = (1 - \eta) \Phi\left(\frac{x-\mu}{\sigma}\right) + \eta G(x), \quad (63)$$

where η is usually interpreted as probability of having an outlying observation and G is a distribution (hopefully symmetric around μ) of outlying observations. It was found that if η is 'small' then using sample median is too pessimistic (and inefficient). We will mention here several alternative options.

Before we proceed note that both the sample mean \bar{X}_n and the sample median \tilde{m}_n can be viewed as M -estimators as

$$\bar{X}_n = \arg \min_{\theta \in \mathbb{R}} \sum_{i=1}^n (X_i - \theta)^2 \quad \text{and} \quad \tilde{m}_n = \arg \min_{\theta \in \mathbb{R}} \sum_{i=1}^n |X_i - \theta|. \quad (64)$$

Huber estimator

This estimator is defined as

$$\hat{\theta}_n^{(H)} = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \rho_H(X_i - \theta),$$

* *Robustní odhad polohy*

where

$$\rho_H(x) = \begin{cases} \frac{x^2}{2}, & |x| \leq k, \\ k \cdot (|x| - \frac{k}{2}), & |x| > k \end{cases} \quad (65)$$

and k is a given constant. Note that the ‘score function’ $\psi_H(x) = \rho'_H(x)$ of the estimator is

$$\psi_H(x) = \rho'_H(x) = \begin{cases} x, & |x| \leq k, \\ k \cdot \text{sign}(x), & |x| > k. \end{cases} \quad (66)$$

Thus one can see that for $x \in (-k, k)$ the function ψ_H corresponds to a score function of a sample mean (which is $\psi(x) = x$) while for $x \in (-\infty, k) \cup (k, \infty)$ it corresponds to a score function of a sample median (which is $\psi(x) = \text{sign}(x)$). Thus Huber estimator presents a compromise between a sample mean and a sample median. So it is not surprising that $\hat{\theta}_n^{(H)}$ is usually a value between the sample median and the sample mean.

When using Huber estimator one has to keep in mind that the identified parameter is

$$\theta^{(H)} = \arg \min_{\theta \in \mathbb{R}} \mathbf{E} \rho_H(X_1 - \theta).$$

Thus if the distribution F is **not symmetric** then $\mathbf{E} X_1$ generally does not coincide with $F^{-1}(0.5)$ and $\theta^{(H)}$ lies between $\mathbf{E} X_1$ and $F^{-1}(0.5)$.

On the other hand if the distribution F is **symmetric**, then $\theta^{(H)}$ coincides with the centre of symmetry, i.e. with $F^{-1}(0.5)$ (the median of F) and also with $\mathbf{E} X_1$, if the expectation exists. It was observed that for the contamination model (63) with G symmetric, Huber estimator usually performs better than the sample mean as well as the sample median. This can be proved analytically by showing that for $\eta > 0$ and G heavy tailed, then usually

$$\text{avar}(\hat{\theta}_n^{(H)}) < \min \{ \text{var}(\bar{X}_n), \text{avar}(\tilde{m}_n) \},$$

where the asymptotic variance $\text{avar}(\hat{\theta}_n^{(H)})$ is derived in Example 42.

The nice thing about Huber estimator is that its loss function $\rho(x; \theta) = \rho_H(x - \theta)$ is convex (in θ) thus $\hat{\theta}_n^{(H)}$ is not too difficult to calculate and with the help of Theorem 10 one can derive its asymptotic distribution (see also Example 42).

Example 42. With the help of Theorem 10 one can show that (under appropriate regularity assumptions)

$$\sqrt{n} (\hat{\theta}_n^{(H)} - \theta^{(H)}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(0, \frac{\sigma_\psi^2}{\gamma^2}),$$

where

$$\gamma = \frac{\partial^2 \mathbf{E} \rho_H(X_1 - \theta)}{\partial \theta^2} \Big|_{\theta = \theta^{(H)}} = F(\theta^{(H)} + k) - F(\theta^{(H)} - k)$$

and

$$\sigma_\psi^2 = \text{var}(\psi_H(X_1 - \theta^{(H)})) = \int_{\theta^{(H)} - k}^{\theta^{(H)} + k} (x - \theta^{(H)})^2 dF(x) + k^2(1 - F(\theta^{(H)} + k) + F(\theta^{(H)} - k)).$$

Thus $\text{avar}(\hat{\theta}_n^{(H)}) = \frac{\sigma_\psi^2}{n\gamma^2}$.

The choice of the constant k is usually done as follows. Suppose that X_1, \dots, X_n follows $N(0, 1)$. Then one takes the smallest k such that

$$\frac{\text{avar}(\hat{\theta}_n^{(H)})}{\text{var}(\bar{X}_n)} \leq 1 + \delta,$$

where δ stands for the efficiency loss of Huber estimator under normal distributions. For instance the common choices are $\delta = 0.05$ or $\delta = 0.1$ which corresponds approximately to $k = 1.37$ or $k = 1.03$.

Other robust M/Z -estimators of location

The other most common M/Z -estimators are the following.

- (i) **Cauchy-pseudolikelihood:** $\rho(x; \theta) = \log(1 + (x - \theta)^2)$. The problem is that this function is not convex in θ and the estimating equation

$$\frac{1}{n} \sum_{i=1}^n \underbrace{\frac{2(X_i - \hat{\theta}_n)}{1 + (X_i - \hat{\theta}_n)^2}}_{\psi(X_i; \hat{\theta}_n)} \stackrel{!}{=} 0$$

has usually more roots.

- (ii) **Tukey's biweight:**

$$\psi(x) = \begin{cases} x \left(1 - \frac{x^2}{k^2}\right)^2, & |x| \leq k, \\ 0, & |x| > k. \end{cases}$$

But also here the corresponding loss function ρ ($\psi = \rho'$) is not convex.

4.2 Robust studentized M/Z -estimators of location

The problem is that the M/Z -estimators presented above (except for the sample mean and the sample median) are not scale equivariant (i.e. $\hat{\theta}_n(c\mathbf{X}) \neq c\hat{\theta}_n(\mathbf{X})$ for each $c \in \mathbb{R}$). That is why in practice M/Z -estimators are usually defined as

$$\hat{\theta}_n = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \rho\left(\frac{X_i - \theta}{S_n}\right), \text{ or as } \sum_{i=1}^n \psi\left(\frac{X_i - \hat{\theta}_n}{S_n}\right) \stackrel{!}{=} 0,$$

where S_n is an appropriate estimator of scale*, which satisfies $S_n(c\mathbf{X}) = |c|S_n(\mathbf{X})$ for each $c \in \mathbb{R}$. The most common estimators of scale are as follows.

* *odhad měřítka*

Sample standard deviation

$$S_n = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2}.$$

Note that in robust statistics S_n is rather rarely used as it is not robust (i.e. it is sensitive to outlying observations).

Interquartile range*

$$S_n = IQR = \hat{F}_n^{-1}(0.75) - \hat{F}_n^{-1}(0.25),$$

where \hat{F}_n is the empirical distribution function (i.e. $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq x\}$). Some people prefer to use

$$\tilde{S}_n = \frac{\hat{F}_n^{-1}(0.75) - \hat{F}_n^{-1}(0.25)}{\Phi^{-1}(0.75) - \Phi^{-1}(0.25)},$$

as it is desired that \tilde{S}_n estimates σ , when X_1, \dots, X_n is a random sample from $\mathbf{N}(\mu, \sigma^2)$.

Note that the breakdown point of interquartile range is 0.25.

Median absolute deviation†

This measure is given as the median absolute deviation from the median, i.e.

$$MAD = \text{med}_{1 \leq i \leq n} \{|X_i - \hat{F}_n^{-1}(0.5)|\},$$

or its modification

$$\widetilde{MAD} = \frac{MAD}{\Phi^{-1}(0.75)},$$

so that it estimates σ for random samples from $\mathbf{N}(\mu, \sigma^2)$.

Note that the breakdown point of this estimator is 0.50.

Remark 16. Note that due to the studentization the functions $\rho(x; \theta) = \rho\left(\frac{x-\theta}{S_n}\right)$ and $\psi(x; \theta) = \psi\left(\frac{x-\theta}{S_n}\right)$ (when viewed as functions of x and θ) are random. Thus one can use neither Theorem 9 nor Theorem 10 to derive the asymptotic distribution of studentized M/Z -estimators.

Nevertheless, if $S_n \xrightarrow[n \rightarrow \infty]{P} S(F)$ and the distribution F is **symmetric**, then (under some regularity assumptions) it can be shown that the asymptotic distribution of studentized Z/M -estimators is the same as the asymptotic distribution of M/Z -estimators with $\rho(x; \theta) = \rho\left(\frac{x-\theta}{S(F)}\right)$ and $\psi(x; \theta) = \psi\left(\frac{x-\theta}{S(F)}\right)$ for which one can (usually) use either Theorem 9 or Theorem 10.

* mezikvartilové rozpětí † mediánová absolutní odchylka

4.3 Robust estimation in linear models

Suppose we observe independent random vectors $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ each of them having the same distribution as the generic random vector (\mathbf{X}, Y) .

4.3.1 The least squares method

This method results in the estimator

$$\hat{\boldsymbol{\beta}}_n^{(LS)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \mathbf{b})^2 = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i Y_i \right).$$

Note that if $\mathbf{X}_{ik} \neq 0$ then by changing Y_i one can arrive at any arbitrary value of $\hat{\beta}_{nk}$.

From Example 40 we know that the method of the least squares identifies the parameter

$$\boldsymbol{\beta}_X^{(LS)} = [\mathbb{E} \mathbf{X} \mathbf{X}^\top]^{-1} \mathbb{E} \mathbf{X} Y$$

and it holds that

$$\sqrt{n} (\hat{\boldsymbol{\beta}}_n^{(LS)} - \boldsymbol{\beta}_X^{(LS)}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(0, \mathbb{V}), \text{ where } \mathbb{V} = [\mathbb{E} \mathbf{X} \mathbf{X}^\top]^{-1} [\mathbb{E} \sigma^2(\mathbf{X}) \mathbf{X} \mathbf{X}^\top] [\mathbb{E} \mathbf{X} \mathbf{X}^\top]^{-1},$$

with $\sigma^2(\mathbf{X}_1) = \mathbb{E} [(Y_1 - \mathbf{X}_1^\top \boldsymbol{\beta}_X)^2 | \mathbf{X}_1]$. Further provided $\mathbb{E} [Y | \mathbf{X}] = \mathbf{X}^\top \boldsymbol{\beta}_0$, then $\boldsymbol{\beta}_X^{(LS)} = \boldsymbol{\beta}_0$ and $\sigma^2(\mathbf{X}) = \text{var}(Y | \mathbf{X})$.

Suppose now that the first component of \mathbf{X}_i is 1 (i.e. the model includes an intercept) and denote by $\widetilde{\mathbf{X}}_i$ the remaining components of \mathbf{X}_i . That is $\mathbf{X}_i = \begin{pmatrix} 1 \\ \widetilde{\mathbf{X}}_i \end{pmatrix}$. Further suppose that the following model holds

$$Y = \beta_0 + \widetilde{\mathbf{X}}^\top \boldsymbol{\beta} + \varepsilon, \text{ where } \varepsilon \perp \widetilde{\mathbf{X}}. \quad (67)$$

Then $\mathbb{E} [Y | \mathbf{X}] = \beta_0 + \boldsymbol{\beta}^\top \widetilde{\mathbf{X}} + \mathbb{E} \varepsilon$ and the method of the least squares identifies the parameter

$$\boldsymbol{\beta}_X^{(LS)} = \begin{pmatrix} \beta_0 + \mathbb{E} \varepsilon \\ \boldsymbol{\beta} \end{pmatrix}. \quad (68)$$

Further the asymptotic variance matrix \mathbb{V} simplifies to

$$\mathbb{V} = \sigma^2 (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1}, \text{ where } \sigma^2 = \text{var}(\varepsilon). \quad (69)$$

4.3.2 Method of the least absolute deviation*

This method is usually considered as a robust alternative to the least squares methods. The estimate of the regression parameter is given by

$$\hat{\boldsymbol{\beta}}_n^{(LAD)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n |Y_i - \mathbf{X}_i^\top \mathbf{b}|,$$

* *Metoda nejmenších absolutních odchylek, mediánová regrese*

As we will see later (see Chapter 5) the LAD method models $\text{med}[Y | \mathbf{X}] = F_{Y|\mathbf{X}}^{-1}(0.5)$ as $\mathbf{X}^\top \boldsymbol{\beta}$. So if indeed $\text{med}[Y | \mathbf{X}] = \mathbf{X}^\top \boldsymbol{\beta}_0$, then $\boldsymbol{\beta}_X^{(LAD)} = \boldsymbol{\beta}_0$.

The *asymptotic distribution* of $\widehat{\boldsymbol{\beta}}_n^{(LAD)}$ can be heuristically derived by Theorem 10 as follows. The score function is given by

$$\psi(\mathbf{x}, y; \mathbf{b}) = -\text{sign}(y - \mathbf{x}^\top \mathbf{b}) \mathbf{x}.$$

Now put $M(\mathbf{b}) = \text{E} [|Y - \mathbf{X}^\top \mathbf{b}| - |Y|]$, where the random vector $\begin{pmatrix} Y \\ \mathbf{X} \end{pmatrix}$ has the same distribution as $\begin{pmatrix} Y_i \\ \mathbf{X}_i \end{pmatrix}$. Then

$$\begin{aligned} \frac{\partial M(\mathbf{b})}{\partial \mathbf{b}} &= \text{E} [\text{sign}(Y - \mathbf{X}^\top \mathbf{b}) (-\mathbf{X})] = -\text{E} \mathbf{X} [\mathbb{1}\{Y > \mathbf{X}^\top \mathbf{b}\} - \mathbb{1}\{Y < \mathbf{X}^\top \mathbf{b}\}] \\ &= -\text{E} \mathbf{X} [1 - 2F_{Y|\mathbf{X}}(\mathbf{X}^\top \mathbf{b})]. \end{aligned}$$

Thus

$$\frac{\partial^2 M(\mathbf{b})}{\partial \mathbf{b} \partial \mathbf{b}^\top} = 2 \text{E} \mathbf{X} f_{Y|\mathbf{X}}(\mathbf{X}^\top \mathbf{b}) \mathbf{X}^\top,$$

which finally implies that

$$\mathcal{J}(\boldsymbol{\beta}_X^{(LAD)}) = \left. \frac{\partial^2 M(\mathbf{b})}{\partial \mathbf{b} \partial \mathbf{b}^\top} \right|_{\mathbf{b}=\boldsymbol{\beta}_X^{(LAD)}} = 2 \text{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(\mathbf{X}^\top \boldsymbol{\beta}_X^{(LAD)})] = 2 \text{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(0.5))].$$

Further as

$$\Sigma(\boldsymbol{\beta}_X^{(LAD)}) = \text{var}(\psi(\mathbf{X}, Y; \boldsymbol{\beta}_X^{(LAD)})) = \text{E} \mathbf{X} \mathbf{X}^\top,$$

one gets that under appropriate regularity assumptions

$$\sqrt{n} (\widehat{\boldsymbol{\beta}}_n^{(LAD)} - \boldsymbol{\beta}_X^{(LAD)}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{V}),$$

where

$$\mathbb{V} = \frac{1}{4} \left(\text{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(0.5))] \right)^{-1} \text{E} \mathbf{X} \mathbf{X}^\top \left(\text{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(0.5))] \right)^{-1}.$$

Note that if model (67) holds, then $\text{med}(Y_1 | \mathbf{X}_1) = \beta_0 + \widetilde{\mathbf{X}}_1^\top \boldsymbol{\beta} + F_\varepsilon^{-1}(0.5)$, where F_ε^{-1} is the quantile function of ε_1 and thus

$$\boldsymbol{\beta}_X^{(LAD)} = \begin{pmatrix} \beta_0 + F_\varepsilon^{-1}(0.5) \\ \boldsymbol{\beta} \end{pmatrix}.$$

Thus when compared with the method of the least squares (68) one can see, that if model (67) holds then both methods identify the same slope parameter $\boldsymbol{\beta}$. The only difference is in intercept.

Further if model (67) holds then

$$f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(0.5)) = f_\varepsilon(F_\varepsilon^{-1}(0.5)),$$

which implies that

$$\mathbb{F}(\boldsymbol{\beta}_X^{(LAD)}) = 2 f_\varepsilon(F_\varepsilon^{-1}(0.5)) \mathbb{E} \mathbf{X} \mathbf{X}^\top$$

and

$$\mathbb{V} = \frac{1}{4[f_\varepsilon(F_\varepsilon^{-1}(0.5))]^2} (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1}. \quad (70)$$

Now when one compares (69) with (70), one can see that the least absolute deviation method is favourable if

$$\frac{1}{[4f_\varepsilon(F_\varepsilon^{-1}(0.5))]^2} < \text{var}(\varepsilon).$$

Regarding the robustness of the least absolute deviation estimator note that in this special situation (i.e. if model (67) holds) $Y_i - \mathbf{X}_i^\top \boldsymbol{\beta}_X^{(LAD)} = \varepsilon_i - F_\varepsilon^{-1}(0.5)$ and the asymptotic representation (50) of $\widehat{\boldsymbol{\beta}}_n$ implies

$$\widehat{\boldsymbol{\beta}}_n^{(LAD)} - \boldsymbol{\beta}_X^{(LAD)} = \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E} \mathbf{X}_1 \mathbf{X}_1^\top \right)^{-1} \mathbf{X}_i \frac{\text{sign}(\varepsilon_i - F_\varepsilon^{-1}(0.5))}{2f_\varepsilon(F_\varepsilon^{-1}(0.5))} + o_P\left(\frac{1}{\sqrt{n}}\right).$$

Thus one can expect that the change of Y_i (or equivalently the change of ε_i) has only a bounded effect on $\widehat{\boldsymbol{\beta}}_n^{(LAD)}$. On the other hand note that the change of \mathbf{X}_i has an unbounded effect on $\widehat{\boldsymbol{\beta}}_n^{(LAD)}$. Thus LAD method is robust with respect to the response but not with respect to the covariates.

The end of
class 17
(19. 4. 2024)

4.3.3 Huber estimator of regression

Analogously as Huber estimator of location is a compromise between a sample mean and a sample median, Huber estimator of regression is a compromise between LS and LAD. Put

$$\widehat{\boldsymbol{\beta}}_n^{(H)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho_H(Y_i - \mathbf{X}_i^\top \mathbf{b}),$$

where ρ_H is defined in (65). Generally, it is difficult to interpret what is being modelled with Huber estimator of regression (it is something between $\mathbb{E}(Y | \mathbf{X})$ and $\text{med}(Y | \mathbf{X})$). Note that it identifies

$$\boldsymbol{\beta}_X^{(H)} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \mathbb{E} \rho_H(Y - \mathbf{X}^\top \mathbf{b}).$$

Equivalently $\boldsymbol{\beta}_X^{(H)}$ solves

$$\mathbb{E} [\psi_H(Y - \mathbf{X}^\top \boldsymbol{\beta}_X^{(H)}) \mathbf{X}] \stackrel{!}{=} \mathbf{0}_p,$$

where ψ_H is defined in (66).

Analogously as in Example 42 one can derive that under appropriate assumptions

$$\sqrt{n} (\widehat{\boldsymbol{\beta}}_n^{(H)} - \boldsymbol{\beta}_X^{(H)}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{V}), \quad \text{with} \quad \mathbb{V} = \mathbb{F}^{-1}(\boldsymbol{\beta}_X^{(H)}) \mathbb{Z}(\boldsymbol{\beta}_X^{(H)}) \mathbb{F}^{-1}(\boldsymbol{\beta}_X^{(H)}),$$

where

$$\mathbb{F}(\boldsymbol{\beta}_X^{(H)}) = \mathbb{E}_{\mathbf{X}} \left[F_{Y|\mathbf{X}}(\mathbf{X}^\top \boldsymbol{\beta}_X^{(H)} + k) - F_{Y|\mathbf{X}}(\mathbf{X}^\top \boldsymbol{\beta}_X^{(H)} - k) \right] \mathbf{X} \mathbf{X}^\top$$

and

$$\mathbb{Z}(\boldsymbol{\beta}_X^{(H)}) = \mathbb{E}_{\mathbf{X}} \left[\mathbf{X} \mathbf{X}^\top \text{var}(\psi(Y - \mathbf{X}^\top \boldsymbol{\beta}_X^{(H)}) | \mathbf{X}) \right].$$

If model (67) holds then $\boldsymbol{\beta}_X^{(H)} = \begin{pmatrix} \beta_{X_0}^{(H)} \\ \boldsymbol{\beta}_X \end{pmatrix}$ solves

$$\mathbb{E} \left[\psi_H(\beta_0 + \widetilde{\mathbf{X}}^\top \boldsymbol{\beta} + \varepsilon - \beta_{X_0}^{(H)} - \widetilde{\mathbf{X}}^\top \widetilde{\boldsymbol{\beta}}_X^{(H)}) \mathbf{X} \right] \stackrel{!}{=} \mathbf{0}_p.$$

Thus $\boldsymbol{\beta}_X$ identifies the following parameter

$$\boldsymbol{\beta}_X^{(H)} = \begin{pmatrix} \beta_0 + \theta^{(H)} \\ \boldsymbol{\beta} \end{pmatrix},$$

where $\theta^{(H)}$ solves $\mathbb{E} \psi_H(\varepsilon - \theta^{(H)}) \stackrel{!}{=} 0$. So if model (67) holds then the interpretation of the regression slope coefficient ($\boldsymbol{\beta}$) is the same for each of the methods described above (LS, LAD, Huber regression).

Further the asymptotic variance matrix simplifies to

$$\mathbb{V} = \frac{\sigma_\psi^2}{\gamma^2} (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1}, \quad (71)$$

with

$$\gamma = F_\varepsilon(\theta^{(H)} + k) - F_\varepsilon(\theta^{(H)} - k)$$

and

$$\sigma_\psi^2 = \int_{\theta^{(H)} - k}^{\theta^{(H)} + k} (x - \theta^{(H)})^2 dF_\varepsilon(x) + k^2 (1 - F_\varepsilon(\theta^{(H)} + k) + F_\varepsilon(\theta^{(H)} - k)).$$

Using (69), (70) and (71) one sees that to compare the efficiency of the estimators $\widehat{\boldsymbol{\beta}}_n^{(LS)}$, $\widehat{\boldsymbol{\beta}}_n^{(LAD)}$ and $\widehat{\boldsymbol{\beta}}_n^{(H)}$ it is sufficient to compare $\text{var}(\varepsilon)$, $\frac{1}{4f_\varepsilon^2(F_\varepsilon^{-1}(0.5))}$ and $\frac{\sigma_\psi^2}{\gamma^2}$.

Regarding the robustness properties the influence function is given by

$$IF(\mathbf{x}, y) = (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1} \frac{1}{\gamma} \psi_H(y - \mathbf{x}^\top \boldsymbol{\beta}_X^{(H)}) \mathbf{x},$$

thus the estimator is robust in response but not in the covariate.

4.3.4 Studentized Huber estimator of regression

Analogously as in Chapter 4.2 in practice the studentized Huber estimator is usually used. This estimator is defined as

$$\widehat{\boldsymbol{\beta}}_n = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho_H \left(\frac{Y_i - \mathbf{X}_i^\top \mathbf{b}}{S_n} \right),$$

where S_n is an estimator of scale of ε_i . For instance one can take MAD or IQR calculated from the residuals from LAD regression $\widehat{\varepsilon}_i = Y_i - \mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n^{(LAD)}$.

Inference:

- With the help of Theorem 10 one can show the asymptotic normality of $\widehat{\boldsymbol{\beta}}_n$ of the (non-Studentized) Huber estimator.
- If model (67) holds, then it can be shown, that the estimate of the scale influences only the asymptotic distribution of the estimate of the intercept and not of the slope.

Literature: Maronna et al. [2006] Chapters 2.1-2.2 and Chapters 4.1-4.4.

5 Quantile regression*

Generally speaking, while the least squares method aims at estimating (modelling) a conditional expectation, quantile regression aims at estimating (modelling) a conditional quantile. This is of interest if the covariate may have different effects on different quantiles of the response.

Applications of the quantile regression can be found in medicine (e.g. constructing reference charts), finance (e.g. estimating value at risk), economics (e.g. wage and income studies, modelling household electricity demand) and environment modelling (e.g. modelling flood height).

5.1 Identification of quantiles

For a given $\tau \in (0, 1)$ consider the following loss function

$$\rho_\tau(x) = \tau x \mathbb{1}\{x > 0\} + (1 - \tau)(-x) \mathbb{1}\{x \leq 0\}.$$

Note that for $x \neq 0$ one gets

$$\psi_\tau(x) = \rho'_\tau(x) = \tau \mathbb{1}\{x > 0\} - (1 - \tau) \mathbb{1}\{x < 0\}.$$

For $x = 0$ put $\psi_\tau(0) = 0$.

Lemma 5. *Let the random variable X have a cumulative distribution function F . Then*

$$F^{-1}(\tau) = \arg \min_{\theta \in \mathbb{R}} \mathbf{E} [\rho_\tau(X - \theta) - \rho_\tau(X)]. \quad (72)$$

* *Kvantilová regrese.*

Proof. Put $M(\theta) = \mathbf{E} [\rho_\tau(X - \theta) - \rho_\tau(X)]$. One can calculate

$$\begin{aligned} M(\theta) &= -\mathbf{E} \int_0^\theta \psi_\tau(X - t) dt = -\int_0^\theta \mathbf{E} \psi_\tau(X - t) dt \\ &= -\int_0^\theta \tau \mathbf{P}(X > t) - (1 - \tau) \mathbf{P}(X < t) dt. \\ &= -\int_0^\theta \tau - \tau F(t) - (1 - \tau)F(t) dt. \\ &= -\tau \theta + \int_0^\theta F(t) dt. \end{aligned}$$

Now for each $\theta < F^{-1}(\tau)$

$$\begin{aligned} M'(\theta_-) &= -\tau + F(\theta_-) \leq -\tau + F(\theta) < 0, \\ M'(\theta_+) &= -\tau + F(\theta_+) = -\tau + F(\theta) < 0. \end{aligned}$$

As the function $M(\theta)$ is continuous, this implies that $M(\theta)$ is decreasing on $(-\infty, F^{-1}(\tau))$.

Analogously for $\theta > F^{-1}(\tau)$

$$\begin{aligned} M'(\theta_-) &= -\tau + F(\theta_-) \geq -\tau + F(F^{-1}(\tau)) \geq 0, \\ M'(\theta_+) &= -\tau + F(\theta_+) \geq -\tau + F(F^{-1}(\tau)) \geq 0. \end{aligned}$$

Thus the function $M(\theta)$ is non-decreasing on $(F^{-1}(\tau), +\infty)$. This further implies that $F^{-1}(\tau)$ is the point of the global minimum of the function $M(\theta)$. \square

Remark 17. Suppose we observe a random sample X_1, \dots, X_n . Let \widehat{F}_n be the corresponding empirical distribution function. Then by

$$\frac{1}{n} \sum_{i=1}^n \rho_\tau(X_i - \theta) = \mathbf{E}_{\widehat{F}_n} \rho_\tau(Z - \theta),$$

where the random variable Z has the distribution given by the empirical distribution function \widehat{F}_n and $\mathbf{E}_{\widehat{F}_n}$ stands for the expectation with respect to this distribution.

Thus by Lemma 5

$$\widehat{F}_n^{-1}(\tau) = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \rho_\tau(X_i - \theta).$$

Note that for $\tau = 0.5$ one gets the characterization of the sample median as in (64).

Further note that from the proof of Lemma 5 it follows that the $\arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \rho_\tau(X_i - \theta)$ is not unique if there exists a root of the function $-\tau + \widehat{F}_n(\theta)$. This happens if $n\tau = i_0 \in \mathbb{N}$ and $X_{(i_0)} < X_{(i_0+1)}$. Then $M(\theta)$ is minimised by any value from the interval $[X_{(i_0)}, X_{(i_0+1)}]$. In this situation $\widehat{F}_n^{-1}(\tau) = X_{(i_0)}$ is the left point of this interval.

5.2 Regression quantiles*

Suppose that one observes independent and identically distributed random vectors

$$\begin{pmatrix} \mathbf{X}_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{X}_n \\ Y_n \end{pmatrix}$$

being distributed as the generic vector $\begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix}$.

The τ -th regression quantile is defined as

$$\widehat{\boldsymbol{\beta}}_n(\tau) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{X}_i^\top \mathbf{b}).$$

At the population level the regression quantile identifies the parameter

$$\boldsymbol{\beta}_X(\tau) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \mathbb{E} [\rho_\tau(Y - \mathbf{X}^\top \mathbf{b}) - \rho_\tau(Y)].$$

Note that thanks to (72)

$$\begin{aligned} \mathbb{E} [\rho_\tau(Y - \mathbf{X}^\top \mathbf{b}) - \rho_\tau(Y)] &= \mathbb{E} \left\{ \mathbb{E} [\rho_\tau(Y - \mathbf{X}^\top \mathbf{b}) - \rho_\tau(Y) \mid \mathbf{X}] \right\} \\ &\geq \mathbb{E} \left\{ \mathbb{E} [\rho_\tau(Y - F_{Y|\mathbf{X}}^{-1}(\tau)) - \rho_\tau(Y) \mid \mathbf{X}] \right\} = \mathbb{E} [\rho_\tau(Y - F_{Y|\mathbf{X}}^{-1}(\tau)) - \rho_\tau(Y)], \end{aligned}$$

where $F_{Y|\mathbf{X}}^{-1}(\tau)$ is the τ -th conditional quantile of Y given \mathbf{X} . Thus if the model for $F_{Y|\mathbf{X}}^{-1}(\tau)$ is correctly specified, that is $F_{Y|\mathbf{X}}^{-1}(\tau) = \mathbf{X}^\top \boldsymbol{\beta}_0$, then $\boldsymbol{\beta}_X(\tau) = \boldsymbol{\beta}_0$.

Often in applications we have $\mathbf{X} = (1, \widetilde{\mathbf{X}}^\top)^\top$ and assume that

$$Y = \beta_0 + \widetilde{\mathbf{X}}^\top \boldsymbol{\beta} + \varepsilon, \quad \text{where } \varepsilon \perp \widetilde{\mathbf{X}}. \quad (73)$$

Then $F_{Y|\mathbf{X}}^{-1}(\tau) = \beta_0 + \widetilde{\mathbf{X}}^\top \boldsymbol{\beta} + F_\varepsilon^{-1}(\tau)$, where $F_\varepsilon^{-1}(\tau)$ is the τ -th quantile of the random error ε . Thus provided model (73) holds

$$\boldsymbol{\beta}_X(\tau) = \begin{pmatrix} \beta_0 + F_\varepsilon^{-1}(\tau) \\ \boldsymbol{\beta} \end{pmatrix}.$$

Thus if model (73) holds, then for $\tau_1 \neq \tau_2$ the regression quantiles $\boldsymbol{\beta}_X(\tau_1)$ and $\boldsymbol{\beta}_X(\tau_2)$ differ only in the intercepts. That is the effect of the covariate is the same for all quantiles of the response. But this is not true in general. In fact the regression quantiles are interesting in situations where the effect of the covariate can be different for different quantiles of the response.

As also illustrated by the following simple examples, the regression quantiles gives us a more detailed idea about the effect of the covariate on the response. This can be of interest on its own or as a check that we do not simplify the situation too much by considering only the effect of the covariate on the conditional expectation.

* *Regresní kvantily*

Example 43. To illustrate consider one-dimensional covariate X_i which is generated from the uniform distribution on the interval $(0, 1)$ and the error term ε_i which has an exponential distribution with mean 1 and which is independent of X_i . Further consider the following two models

- *The homoscedastic model* given by

$$Y_i = 1 + 2 X_i + \varepsilon_i, \quad i = 1, \dots, n.$$

- *The heteroscedastic model* given by

$$Y_i = 1 + 2 X_i + 2 X_i \varepsilon_i, \quad i = 1, \dots, n.$$

On Figure 43 one can find a random sample of size 1 000 from these models. The solid lines represent the fitted regression quantiles for $\tau \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$ assuming that the conditional quantile is in the simple linear form

$$F_{Y|X}^{-1}(\tau) = \beta_1(\tau) + \beta_2(\tau) X.$$

The standard least square estimator is included for the reason of comparison.

Note that in the homoscedastic model all the fitted lines are approximately parallel. This is in agreement with the above finding that in the ‘strict linear model’ (73) the slope of the (theoretical) regression quantiles is the same (up to the random variations that decreases as the sample size increases).

On the other hand in the heteroscedastic model the slopes differ and in this simple example we see that the effect of the covariate is stronger on larger conditional quantiles.

Homework exercise. In the homoscedastic as well as heteroscedastic model find the theoretical conditional quantile $F_{Y|X}^{-1}(\tau)$ for different values of τ and compare it with the conditional expectation $E[Y|X]$. Compare the results with the fitted lines on Figure 43.

Example 44. Let Y_1, \dots, Y_{n_1} be a random sample with the distribution function F and $Y_{n_1+1}, \dots, Y_{n_1+n_2}$ be a random sample from the distribution function G .

Often it is assumed that $G(x) = F(x + \mu)$ for each $x \in \mathbb{R}$. Thus alternatively we can formulate the two-sample problem as a linear regression problem with

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \tag{74}$$

where

$$x_i = \begin{cases} 0, & i = 1, \dots, n_1, \\ 1, & i = n_1 + 1, \dots, n_1 + n_2, \end{cases}$$

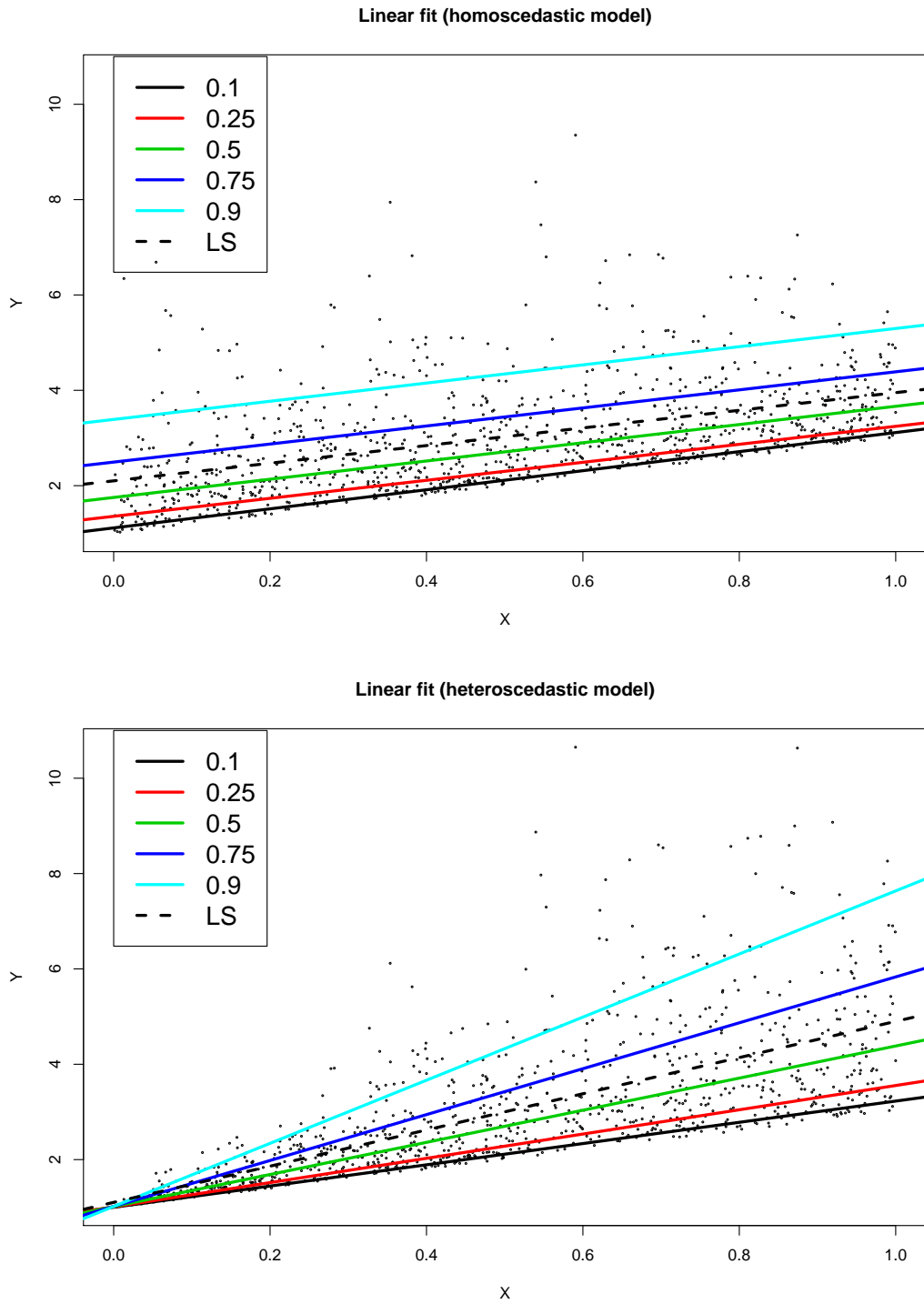


Figure 1: Fitted regression quantiles for $\tau \in \{0.1, 0.25, 0.5, 0.75, 0.9\}$ (solid lines with different colours) for homoscedastic model (the upper figure) and heteroscedastic model (the lower figure). The least squares fit is included for the reason of comparison (dashed line).

and the error term ε_i has a cumulative distribution function F . Usually we are interested in estimating β_1 . By the LS method one gets

$$\widehat{\beta}_1 = \frac{1}{n_2} \sum_{i=n_1+1}^{n_1+n_2} Y_i - \frac{1}{n_1} \sum_{i=1}^{n_1} Y_i \xrightarrow[n_1, n_2 \rightarrow \infty]{P} \underbrace{\mu_G - \mu_F}_{=:\mu} =: \beta_1^{LS},$$

where μ_F and μ_G stand for the expectation of an observation from the first and second sample respectively.

On the other hand let $n = n_1 + n_2$. Then the quantile regression yields

$$\begin{aligned} \widehat{\beta}(\tau) &= \arg \min_{b_0, b_1} \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - b_0 - b_1 x_i) \\ &= \arg \min_{b_0, b_1} \frac{1}{n} \left(\sum_{i=1}^{n_1} \rho_\tau(Y_i - b_0) + \sum_{i=n_1+1}^{n_1+n_2} \rho_\tau(Y_i - b_0 - b_1) \right). \end{aligned}$$

The first sum is minimised by

$$\widehat{\beta}_0(\tau) = F_{n_1}^{-1}(\tau)$$

and the second sum by

$$\beta_0(\tau) + \widehat{\beta}_1(\tau) = G_{n_2}^{-1}(\tau).$$

Thus we get

$$\widehat{\beta}_1(\tau) = G_{n_2}^{-1}(\tau) - F_{n_1}^{-1}(\tau) \xrightarrow[n_1, n_2 \rightarrow \infty]{P} G^{-1}(\tau) - F^{-1}(\tau) := \beta_1(\tau).$$

Further if model (74) really holds, then $G^{-1}(\tau) = F^{-1}(\tau) + \mu$ and one gets $\beta_1(\tau) = \mu = \beta_1^{LS}$ for each $\tau \in (0, 1)$.

Computing regression quantiles

The optimisation task

$$\min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho_\tau(Y_i - \mathbf{X}_i^T \mathbf{b})$$

can be rewritten with the help of linear programming as minimisation of the objective function

$$\tau \sum_{i=1}^n r_i^+ + (1 - \tau) \sum_{i=1}^n r_i^-,$$

subject to the following constrains

$$\begin{aligned} \sum_{j=1}^p X_{ij} b_j + r_i^+ - r_i^- &= Y_i, & i = 1, \dots, n, \\ r_i^+ &\geq 0, \quad r_i^- \geq 0, & i = 1, \dots, n, \\ b_j &\in \mathbb{R}, & j = 1, \dots, p. \end{aligned}$$

Note that one can think of r_i^+ and r_i^- as the positive or negative part of the i -th residual, i.e.

$$r_i^+ = (Y_i - \mathbf{X}_i^\top \mathbf{b})_+, \quad r_i^- = (Y_i - \mathbf{X}_i^\top \mathbf{b})_-.$$

This can be solved for instance with the help of *the simplex algorithm*.

5.3 Interpretation of the regression quantiles

Provided $F_{Y|\mathbf{X}}^{-1}(\tau) = \mathbf{X}^\top \boldsymbol{\beta}$ and the model is correctly specified then one can interpret $\hat{\beta}_{nk}(\tau)$ (the k -th element of $\hat{\boldsymbol{\beta}}_n(\tau)$) as the estimated change of the conditional quantile of the response when the k -th element of the explanation variable increases by 1.

Intersection of the fitted regression quantiles

Note that it might happen that for a given value of the covariate \mathbf{x} and given quantiles $0 < \tau_1 < \tau_2 < 1$

$$\hat{F}_{Y|\mathbf{X}=\mathbf{x}}^{-1}(\tau_1) = \mathbf{x}^\top \hat{\boldsymbol{\beta}}_n(\tau_1) > \mathbf{x}^\top \hat{\boldsymbol{\beta}}_n(\tau_2) = \hat{F}_{Y|\mathbf{X}=\mathbf{x}}^{-1}(\tau_2). \quad (75)$$

which is rather strange as we know that the theoretical quantiles for $\tau_1 < \tau_2$ must satisfy

$$F_{Y|\mathbf{X}=\mathbf{x}}^{-1}(\tau_1) \leq F_{Y|\mathbf{X}=\mathbf{x}}^{-1}(\tau_2).$$

Thus if one gets the inequality (75) (we also say that the regression quantiles cross) for \mathbf{x} from the support of the covariate, it might indicate that the assumed linear model for the conditional quantile is not correct.

Transformed response

It is worth noting that if one models the conditional quantile of the transformed response, that is one assumes that $F_{h(Y)|\mathbf{X}}^{-1}(\tau) = \mathbf{X}^\top \boldsymbol{\beta}$ for a given increasing transformation h , then

$$\tau = \mathbb{P}(h(Y) \leq \mathbf{X}^\top \boldsymbol{\beta} | \mathbf{X}) = \mathbb{P}(Y \leq h^{-1}(\mathbf{X}^\top \boldsymbol{\beta}) | \mathbf{X}),$$

which implies that $F_{Y|\mathbf{X}}^{-1}(\tau) = h^{-1}(\mathbf{X}^\top \boldsymbol{\beta})$. Analogously $F_{Y|\mathbf{X}}^{-1}(1 - \tau) = h^{-1}(\mathbf{X}^\top \boldsymbol{\beta})$ for h decreasing. That is unlike for modelling of conditional expectation (through the least squares method), here we still have a link between $\boldsymbol{\beta}$ and the quantile of the original (not transformed) response $F_{Y|\mathbf{X}}^{-1}(\tau)$.

Thus from the practical point of view even if $\hat{\boldsymbol{\beta}}_n(\tau)$ is estimated from the response-transformed data $(\mathbf{X}_1)_{h(Y_1)}, \dots, (\mathbf{X}_n)_{h(Y_n)}$, one can still estimate the conditional quantile of the original (not transformed) data $\hat{F}_{Y|\mathbf{X}}^{-1}(\tau) = h^{-1}(\mathbf{X}^\top \hat{\boldsymbol{\beta}}_n(\tau))$ (for h increasing). On the other hand if we estimate the conditional expectation of $\mathbb{E}[h(Y)|\mathbf{X}]$ as $\mathbf{X}^\top \hat{\boldsymbol{\beta}}_n$, there is no general way how to use $\hat{\boldsymbol{\beta}}_n$ to get an estimate of $\mathbb{E}[Y|\mathbf{X}]$.

A very common and popular transformation is log-transformation, i.e. $h(y) = \log y$. This results in $F_{Y|\mathbf{X}}^{-1}(\tau) = e^{\mathbf{X}^\top \boldsymbol{\beta}(\tau)}$ and $e^{\beta_k(\tau)}$ measures how many times the conditional quantile $F_{Y|\mathbf{X}}^{-1}(\tau)$ changes when the k -th coordinate of the covariate is increased by adding one.

The end of
class 19
(26. 4. 2024)

5.4 Inference for regression quantiles

Analogously as in Chapter 4.3.2 one can heuristically derive that under appropriate regularity assumption for fixed $\tau \in (0, 1)$

$$\sqrt{n} (\widehat{\boldsymbol{\beta}}_n(\tau) - \boldsymbol{\beta}_X(\tau)) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}_p(\mathbf{0}_p, \mathbb{V}),$$

where

$$\mathbb{V} = \left(\mathbb{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(\tau))] \right)^{-1} \tau(1-\tau) \mathbb{E} \mathbf{X} \mathbf{X}^\top \left(\mathbb{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(\tau))] \right)^{-1}. \quad (76)$$

Note that if model (73) holds, then $F_{Y|\mathbf{X}}^{-1}(\tau) = \beta_0 + \widetilde{\mathbf{X}}^\top \boldsymbol{\beta} + F_\varepsilon^{-1}(\tau)$, where F_ε^{-1} is the quantile function of ε_1 and thus

$$\boldsymbol{\beta}_X(\tau) = \begin{pmatrix} \beta_0 + F_\varepsilon^{-1}(\tau) \\ \boldsymbol{\beta} \end{pmatrix}.$$

Further if model (73) holds then

$$f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(\tau)) = f_\varepsilon(F_\varepsilon^{-1}(\tau)),$$

which implies that

$$\mathbb{E} [\mathbf{X} \mathbf{X}^\top f_{Y|\mathbf{X}}(F_{Y|\mathbf{X}}^{-1}(\tau))] = f_\varepsilon(F_\varepsilon^{-1}(\tau)) \mathbb{E} \mathbf{X} \mathbf{X}^\top$$

and

$$\mathbb{V} = \frac{\tau(1-\tau)}{[f_\varepsilon(F_\varepsilon^{-1}(\tau))]^2} (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1}. \quad (77)$$

Estimation of asymptotic variance of $\widehat{\boldsymbol{\beta}}_n(\tau)$

Note that in general the asymptotic variance matrix (76) of $\widehat{\boldsymbol{\beta}}_n(\tau)$ is rather complicated and it is not clear how to estimate it. That is why nonparametric bootstrap is of interest.

If model (73) holds, then the asymptotic variance matrix of $\widehat{\boldsymbol{\beta}}_n(\tau)$ simplifies considerably and one gets

$$\text{avar}(\widehat{\boldsymbol{\beta}}_n(\tau)) = \frac{1}{n} (\mathbb{E} \mathbf{X} \mathbf{X}^\top)^{-1} \frac{\tau(1-\tau)}{f_\varepsilon^2(F_\varepsilon^{-1}(\tau))}.$$

The matrix $\mathbf{E} \mathbf{X} \mathbf{X}^\top$ can be estimated as $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top$. The difficulty is in estimating the sparsity function $s(\tau) = \frac{1}{f_\varepsilon(F_\varepsilon^{-1}(\tau))}$. In Chapter 4.10.1 of [Koenker \[2005\]](#) it is suggested that one can use the following estimate

$$\widehat{s}_n(\tau) = \frac{\widehat{F}_{n\widehat{\varepsilon}}^{-1}(\tau + h_n) - \widehat{F}_{n\widehat{\varepsilon}}^{-1}(\tau - h_n)}{2 h_n},$$

where

$$\widehat{F}_{n\widehat{\varepsilon}}(y) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{Y_i - \mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n(\tau) \leq y\} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\widehat{\varepsilon}_i(\tau) \leq y\}$$

is the empirical distribution function of the residuals and (the bandwidth) h_n is a sequence going to zero as $n \rightarrow \infty$. A possible choice of h_n (derived when assuming normal errors $\varepsilon_1, \dots, \varepsilon_n$) is given by

$$h_n = n^{-1/3} u_{1-\alpha/2}^{2/3} \left[\frac{1.5 \varphi^2(u_\tau)}{2 u_\tau^2 + 1} \right]^{1/3},$$

where φ is the density of $\mathbf{N}(0, 1)$. For details and other possible choices of h_n see Chapter 4.10.1 in [Koenker \[2005\]](#) and the references therein.

As estimating $\frac{1}{f_\varepsilon(F_\varepsilon^{-1}(\tau))}$ is rather delicate, also in this situation the nonparametric bootstrap* is of interest.

5.5 Asymptotic normality of sample quantiles†

Suppose that we have a random sample X_1, \dots, X_n , where X_1 has a cumulative distribution function F . Note that for a given $\tau \in (0, 1)$ thanks to Remark 17 one can view the sample quantile $\widehat{F}_n^{-1}(\tau)$ as the argument of minimum of a convex function. Thus analogously as in Chapter 3.4 one can derive that if $f(x)$ (the density of X_1) is positive and continuous in a neighbourhood of $F^{-1}(\tau)$, then

$$\sqrt{n} (\widehat{F}_n^{-1}(\tau) - F^{-1}(\tau)) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}\left(0, \frac{\tau(1-\tau)}{f^2(F^{-1}(\tau))}\right).$$

Literature: [Koenker \[2005\]](#), Sections 2.1, 2.4, 4.2, 4.10.

6 EM-algorithm

It is an *iterative* algorithm to find the *maximum likelihood* estimator $\widehat{\boldsymbol{\theta}}_n$ in situations with missing data. It is also often used in situations when the model can be specified with the help of some unobserved variables and finding $\widehat{\boldsymbol{\theta}}_n$ would be (relatively) simple with the knowledge of those unobserved variables.

* Bootstrap and other resampling methods are in detail explained in the course [NMST545 Mathematical Statistics 4](#). † Not done at the lecture. It is assumed that it is known from the bachelor degree.

Example 45. Let X_1, \dots, X_n be a random sample from the distribution with the density

$$f(x; \boldsymbol{\pi}) = \sum_{j=1}^G \pi_j f_j(x),$$

where f_1, \dots, f_G are known densities and $\boldsymbol{\pi} = (\pi_1, \dots, \pi_G)^\top$ is a vector of unknown non-negative *mixing proportions* such that $\sum_{j=1}^G \pi_j = 1$. Find the maximum likelihood estimator of the parameter $\boldsymbol{\pi}$, i.e.

$$\hat{\boldsymbol{\pi}}_n = \arg \max_{\boldsymbol{\pi} \in \Theta} \left(\prod_{i=1}^n f(X_i; \boldsymbol{\pi}) \right),$$

where $\Theta = \{(\pi_1, \dots, \pi_G)^\top : \pi_j \in [0, 1], \sum_{j=1}^G \pi_j = 1\}$.

Solution. A straightforward approach would be to maximize the log-likelihood

$$\ell_n(\boldsymbol{\pi}) = \sum_{i=1}^n \log f(X_i; \boldsymbol{\pi}) = \sum_{i=1}^n \log \left(\sum_{j=1}^G \pi_j f_j(X_i) \right).$$

Using for instance the parametrization $\pi_G = 1 - \sum_{j=1}^{G-1} \pi_j$, the system of score equations is given by

$$U_{jn}(\boldsymbol{\pi}) = \frac{\partial \ell_n(\boldsymbol{\pi})}{\partial \pi_j} = \sum_{i=1}^n \left[\frac{f_j(X_i)}{\sum_{l=1}^G \pi_l f_l(X_i)} - \frac{f_G(X_i)}{\sum_{l=1}^G \pi_l f_l(X_i)} \right] \stackrel{!}{=} 0, \quad j = 1, \dots, G-1,$$

which requires some numerical routines.

Alternatively one can use the EM-algorithm, which runs as follows. Introduce $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{iG})^\top \sim \text{Mult}_G(1; \boldsymbol{\pi})$, where

$$Z_{ij} = \begin{cases} 1, & X_i \text{ is generated from } f_j(x), \\ 0, & \text{otherwise.} \end{cases}$$

Note that one can think of our data as the realizations of the independent and identically distributed random vectors $\begin{pmatrix} X_1 \\ \mathbf{Z}_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ \mathbf{Z}_n \end{pmatrix}$, where $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ are missing.

Put $\mathbb{X} = (X_1, \dots, X_n)^\top$. The joint density of a random vector $\begin{pmatrix} X_i \\ \mathbf{Z}_i \end{pmatrix}$ is given by

$$f_{X, \mathbf{Z}}(x, \mathbf{z}; \boldsymbol{\pi}) = f_{X|\mathbf{Z}}(x|\mathbf{z}; \boldsymbol{\pi}) f_{\mathbf{Z}}(\mathbf{z}; \boldsymbol{\pi}) = \left(\sum_{j=1}^G z_j f_j(x) \right) \cdot \left(\prod_{j=1}^G \pi_j^{z_j} \right).$$

In the context of EM algorithm the random sample $\begin{pmatrix} X_1 \\ \mathbf{Z}_1 \end{pmatrix}, \dots, \begin{pmatrix} X_n \\ \mathbf{Z}_n \end{pmatrix}$ is called *complete data*.

The corresponding log-likelihood is called *complete log-likelihood* and it is given by

$$\begin{aligned}\ell_n^C(\boldsymbol{\pi}) &= \log \left\{ \prod_{i=1}^n \left[\left(\sum_{j=1}^G Z_{ij} f_j(X_i) \right) \left(\prod_{j=1}^G \pi_j^{Z_{ij}} \right) \right] \right\} \\ &= \sum_{i=1}^n \left[\log \left(\sum_{j=1}^G Z_{ij} f_j(X_i) \right) \right] + \sum_{i=1}^n \left[\sum_{j=1}^G Z_{ij} \log \pi_j \right].\end{aligned}$$

If we knew $\mathbf{Z}_1, \dots, \mathbf{Z}_n$, then we would estimate simply $\hat{\pi}_j = \frac{1}{n} \sum_{i=1}^n Z_{ij}, j = 1, \dots, G$. The EM algorithm runs in the following two steps:

- (i) **E-step** (Expectation step): Let $\hat{\boldsymbol{\pi}}^{(k)}$ be the current estimate of $\boldsymbol{\pi}$. In this step we calculate

$$Q(\boldsymbol{\pi}, \hat{\boldsymbol{\pi}}^{(k)}) = \mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}}[\ell_n^C(\boldsymbol{\pi}) | \mathbb{X}],$$

where the expectation is taken with respect to the unobserved random vectors $\mathbf{Z}_1, \dots, \mathbf{Z}_n$. More precisely one has to take the expectation with respect to the conditional distribution of $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ given X_1, \dots, X_n . As this distribution depends on the unknown parameter $\boldsymbol{\pi}$, this parameter is replaced with the current version of the estimate $\hat{\boldsymbol{\pi}}^{(k)}$. This is indicated by $\mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}}$. Note that in this step one gets rid of the unobserved random vectors $\mathbf{Z}_1, \dots, \mathbf{Z}_n$.

- (ii) **M-step** (Maximization step): The updated value of the estimate of $\boldsymbol{\pi}$ is calculated as

$$\hat{\boldsymbol{\pi}}^{(k+1)} = \arg \max_{\boldsymbol{\pi} \in \Theta} Q(\boldsymbol{\pi}, \hat{\boldsymbol{\pi}}^{(k)}).$$

E-step in a detail:

$$Q(\boldsymbol{\pi}, \hat{\boldsymbol{\pi}}^{(k)}) = \mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}} \left[\sum_{i=1}^n \log \left(\sum_{j=1}^G Z_{ij} f_j(X_i) \right) \middle| \mathbb{X} \right] + \mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}} \left[\sum_{i=1}^n \sum_{j=1}^G Z_{ij} \log \pi_j \middle| \mathbb{X} \right]. \quad (78)$$

Note that the first term on the right-hand side of the above equation does not depend on $\boldsymbol{\pi}$. Thus we do not need to calculate this term for M-step. To calculate the second term it is sufficient to calculate $\mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}} [Z_{ij} | \mathbb{X}]$. To do that denote $\mathbf{e}_j = (0, \dots, 0, 1, 0, \dots, 0)^\top$ for the j -th canonical vector. Now with the help of Bayes theorem for densities (Theorem A15) one can calculate

$$\begin{aligned}\mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}} [Z_{ij} | \mathbb{X}] &= \mathbb{E}_{\hat{\boldsymbol{\pi}}^{(k)}} [Z_{ij} | X_i] = \mathbb{P}_{\hat{\boldsymbol{\pi}}^{(k)}}(Z_{ij} = 1 | X_i) = f_{\mathbf{Z}|X}(\mathbf{e}_j | X_i; \hat{\boldsymbol{\pi}}^{(k)}) \\ &= \frac{f_{X|\mathbf{Z}}(X_i | \mathbf{e}_j; \hat{\boldsymbol{\pi}}^{(k)}) f_{\mathbf{Z}}(\mathbf{e}_j; \hat{\boldsymbol{\pi}}^{(k)})}{f_X(X_i; \hat{\boldsymbol{\pi}}^{(k)})} = \frac{f_j(X_i) \hat{\pi}_j^{(k)}}{\sum_{l=1}^G f_l(X_i) \hat{\pi}_l^{(k)}} =: z_{ij}^{(k)}.\end{aligned}$$

The end of
class 20
(30. 4. 2024)

M-step in a detail: Note that with the help of the previous step and (78)

$$Q(\boldsymbol{\pi}, \widehat{\boldsymbol{\pi}}^{(k)}) = \text{const} + \sum_{i=1}^n \sum_{j=1}^G z_{ij}^{(k)} \log \pi_j.$$

Analogously as when calculating the maximum likelihood estimator in a multinomial distribution one can show that the updated value of the estimate of $\boldsymbol{\pi}$ is given by

$$\widehat{\boldsymbol{\pi}}^{(k+1)} = \arg \max_{\boldsymbol{\pi} \in \Theta} Q(\boldsymbol{\pi}, \widehat{\boldsymbol{\pi}}^{(k)}) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^{(k)},$$

where $\mathbf{z}_i^{(k)} = (z_{i1}^{(k)}, \dots, z_{iG}^{(k)})^\top$ and so $\widehat{\pi}_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n z_{ij}^{(k)}$ for $j \in \{1, \dots, G\}$.

6.1 General description of the EM-algorithm

Denote the observed random variables as \mathbb{Y}_{obs} and the unobserved (missing) random variables \mathbb{Y}_{mis} . Let $f(\mathbf{y}; \boldsymbol{\theta})$ be the joint density (with respect to a σ -finite measure μ) of $\mathbb{Y} = (\mathbb{Y}_{obs}, \mathbb{Y}_{mis})$ and denote $\ell_n^C(\boldsymbol{\theta})$ the *complete log-likelihood* of \mathbb{Y} . Our task is to maximize the *observed log-likelihood* $\ell_{obs}(\boldsymbol{\theta}) = \log f(\mathbb{Y}_{obs}; \boldsymbol{\theta})$, where $f(\mathbf{y}_{obs}; \boldsymbol{\theta})$ is the density of \mathbb{Y}_{obs} . Note that

$$\begin{aligned} \ell_n^C(\boldsymbol{\theta}) &= \log f(\mathbb{Y}_{obs}, \mathbb{Y}_{mis}; \boldsymbol{\theta}) = \log (f(\mathbb{Y}_{mis} | \mathbb{Y}_{obs}; \boldsymbol{\theta}) f(\mathbb{Y}_{obs}; \boldsymbol{\theta})) \\ &= \log f(\mathbb{Y}_{mis} | \mathbb{Y}_{obs}; \boldsymbol{\theta}) + \log f(\mathbb{Y}_{obs}; \boldsymbol{\theta}) = \log f(\mathbb{Y}_{mis} | \mathbb{Y}_{obs}; \boldsymbol{\theta}) + \ell_{obs}(\boldsymbol{\theta}), \end{aligned}$$

where $f(\mathbf{y}_{mis} | \mathbf{y}_{obs}; \boldsymbol{\theta})$ stands for the conditional density of \mathbb{Y}_{mis} given $\mathbb{Y}_{obs} = \mathbf{y}_{obs}$. Thus one can express the observed log-likelihood with the help of the complete log-likelihood as

$$\ell_{obs}(\boldsymbol{\theta}) = \ell_n^C(\boldsymbol{\theta}) - \log f(\mathbb{Y}_{mis} | \mathbb{Y}_{obs}; \boldsymbol{\theta}). \quad (79)$$

Finally denote

$$Q(\boldsymbol{\theta}, \widetilde{\boldsymbol{\theta}}) = \mathbf{E}_{\widetilde{\boldsymbol{\theta}}} [\ell_n^C(\boldsymbol{\theta}) | \mathbb{Y}_{obs}]. \quad (80)$$

EM-algorithm runs as follows:

Let $\widehat{\boldsymbol{\theta}}^{(k)}$ be the result of the k -th iteration of the EM-algorithm. The next iteration $\widehat{\boldsymbol{\theta}}^{(k+1)}$ is computed in two steps:

E-step: Calculate $Q(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)})$.

M-step: Find $\widehat{\boldsymbol{\theta}}^{(k+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)})$.

Note that at this moment it is not at all clear, if the EM-algorithm is a good idea. Remember that our task is to maximize the observed likelihood $\ell_{obs}(\boldsymbol{\theta})$. The following theorem is the first answer in this aspect.

Theorem 11. Let the set $\{y_{miss} : f(y_{miss}|\mathbb{Y}_{obs}; \boldsymbol{\theta}) > 0\}$ does not depend on $\boldsymbol{\theta}$. Further $\ell_{obs}(\boldsymbol{\theta})$ be the observed likelihood and $\widehat{\boldsymbol{\theta}}^{(k)}$ be the result of the k -th iteration of the EM-algorithm. Then

$$\ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k+1)}) \geq \ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k)}).$$

Proof. Note that the left-hand side of (79) does not depend on \mathbb{Y}_{mis} . Thus applying $\mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}[\cdot | \mathbb{Y}_{obs}]$ on both sides of (79) yields that

$$\begin{aligned} \ell_{obs}(\boldsymbol{\theta}) &= \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}[\ell_n^C(\boldsymbol{\theta}) | \mathbb{Y}_{obs}] - \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}[\log f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \boldsymbol{\theta}) | \mathbb{Y}_{obs}] \\ &=: Q(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}) - H(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}). \end{aligned} \quad (81)$$

Now note that

$$\begin{aligned} \ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k+1)}) &= Q(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}) - H(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}), \\ \ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k)}) &= Q(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}) - H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}). \end{aligned}$$

Thus to verify $\ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k+1)}) \geq \ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k)})$ it is sufficient to show that

$$Q(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}) \geq Q(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}) \quad \text{and also} \quad H(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}) \leq H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}). \quad (82)$$

Showing *the first inequality* in (82) is easy as from the M-step

$$\widehat{\boldsymbol{\theta}}^{(k+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}),$$

which implies that $Q(\widehat{\boldsymbol{\theta}}^{(k+1)}, \widehat{\boldsymbol{\theta}}^{(k)}) \geq Q(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)})$ for each $\boldsymbol{\theta} \in \Theta$.

To show *the second inequality* in (82) one gets with the help of Jensen's inequality that for each $\boldsymbol{\theta} \in \Theta$:

$$\begin{aligned} H(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}) &= \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}[\log f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \boldsymbol{\theta}) | \mathbb{Y}_{obs}] \\ &= \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}\left[\log\left(\frac{f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \boldsymbol{\theta})}{f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \widehat{\boldsymbol{\theta}}^{(k)})}\right) \middle| \mathbb{Y}_{obs}\right] + \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}[\log f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \widehat{\boldsymbol{\theta}}^{(k)}) | \mathbb{Y}_{obs}] \\ &\stackrel{\text{Jensen}}{\leq} \log\left(\mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}}\left[\frac{f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \boldsymbol{\theta})}{f(\mathbb{Y}_{mis}|\mathbb{Y}_{obs}; \widehat{\boldsymbol{\theta}}^{(k)})} \middle| \mathbb{Y}_{obs}\right]\right) + H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}) \\ &= \log\left(\int \frac{f(\mathbf{y}_{mis}|\mathbb{Y}_{obs}; \boldsymbol{\theta})}{f(\mathbf{y}_{mis}|\mathbb{Y}_{obs}; \widehat{\boldsymbol{\theta}}^{(k)})} \cdot f(\mathbf{y}_{mis}|\mathbb{Y}_{obs}; \widehat{\boldsymbol{\theta}}^{(k)}) d\mu(\mathbf{y}_{mis})\right) + H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}) \\ &= \log(1) + H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}) = H(\widehat{\boldsymbol{\theta}}^{(k)}, \widehat{\boldsymbol{\theta}}^{(k)}). \end{aligned} \quad (83)$$

□

6.2 Convergence of the EM-algorithm

Although from Theorem 11 we know that EM algorithm increases (more precisely does not decrease) the observed log-likelihood, it is still not clear whether the sequence $\{\widehat{\boldsymbol{\theta}}^{(k)}\}_{k=1}^{\infty}$ converges. And if it converges what is the limit.

To answer this question we need to introduce the following regularity assumptions.

- The parameter space Θ is a subset of \mathbb{R}^p .
- The set $\Theta_0 = \{\boldsymbol{\theta} \in \Theta : \ell_{obs}(\boldsymbol{\theta}) \geq \ell_{obs}(\boldsymbol{\theta}_0)\}$ is compact for any $\boldsymbol{\theta}_0 \in \Theta$ such that $\ell_{obs}(\boldsymbol{\theta}_0) > -\infty$.
- $\ell_{obs}(\boldsymbol{\theta})$ is continuous in Θ and differentiable in the interior of Θ .

Theorem 12. *Let the function $Q(\boldsymbol{\theta}, \widetilde{\boldsymbol{\theta}})$ defined in (80) be continuous both in $\boldsymbol{\theta}$ and $\widetilde{\boldsymbol{\theta}}$. Then all the limit points of any instance $\{\widehat{\boldsymbol{\theta}}^{(k)}\}$ are stationary points of $\ell_{obs}(\boldsymbol{\theta})$. Further $\{\ell_{obs}(\widehat{\boldsymbol{\theta}}^{(k)})\}$ converges monotonically to some value $\ell^* = \ell_{obs}(\boldsymbol{\theta}^*)$, where $\boldsymbol{\theta}^*$ is a stationary point of $\ell_{obs}(\boldsymbol{\theta})$.*

Proof. See Wu [1983]. □

Note that if $\boldsymbol{\theta}^*$ is a stationary point of $\ell_{obs}(\boldsymbol{\theta})$, then

$$\left. \frac{\partial \ell_{obs}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} = \mathbf{0}_p.$$

Thus by Theorem 12 the EM-algorithm finds a solution of the system of log-likelihood equations but in generally there is no guarantee that this is a global maximum of $\ell_{obs}(\boldsymbol{\theta})$.

Corollary 2. *Let the assumptions of Theorem 12 be satisfied. Further suppose that the function $\ell_{obs}(\boldsymbol{\theta})$ has a unique maximum $\widehat{\boldsymbol{\theta}}_n$ that is the only stationary point. Then $\widehat{\boldsymbol{\theta}}^{(k)} \rightarrow \widehat{\boldsymbol{\theta}}_n$ as $k \rightarrow \infty$.*

Example 46. Let X_1, \dots, X_n be a random sample from the distribution with the density

$$f(x) = w \frac{1}{\sigma_1} \varphi\left(\frac{x-\mu_1}{\sigma_1}\right) + (1-w) \frac{1}{\sigma_2} \varphi\left(\frac{x-\mu_2}{\sigma_2}\right),$$

where $w \in [0, 1]$, $\mu_1, \mu_2 \in \mathbb{R}$, $\sigma_1^2, \sigma_2^2 \in (0, \infty)$ are unknown parameters and

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp\{-x^2/2\}$$

is the density of the standard normal distribution. Describe the EM algorithm to find the maximum likelihood estimates of the unknown parameters.

6.3 Rate of convergence of EM-algorithm

Note that in the M-step of the algorithm there might not be a unique value that maximizes $Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)})$. Thus denote the set of maximizing points as $\mathcal{M}(\hat{\boldsymbol{\theta}}^{(k)})$, i.e.

$$\mathcal{M}(\hat{\boldsymbol{\theta}}^{(k)}) = \left\{ \tilde{\boldsymbol{\theta}} : Q(\tilde{\boldsymbol{\theta}}, \hat{\boldsymbol{\theta}}^{(k)}) = \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)}) \right\}.$$

Then one needs to choose $\hat{\boldsymbol{\theta}}^{(k+1)}$ as an element of the set $\mathcal{M}(\hat{\boldsymbol{\theta}}^{(k)})$. Thus let $\mathbf{M} : \Theta \rightarrow \Theta$ be a mapping such that

$$\hat{\boldsymbol{\theta}}^{(k+1)} = \mathbf{M}(\hat{\boldsymbol{\theta}}^{(k)}).$$

Let $\hat{\boldsymbol{\theta}}^{(k)} \rightarrow \boldsymbol{\theta}^*$ as $k \rightarrow \infty$. Note that then $\boldsymbol{\theta}^* = \mathbf{M}(\boldsymbol{\theta}^*)$. Assuming that \mathbf{M} is sufficiently smooth one gets by the one term Taylor expansion around the point $\boldsymbol{\theta}^*$ the following approximation

$$\hat{\boldsymbol{\theta}}^{(k+1)} = \mathbf{M}(\hat{\boldsymbol{\theta}}^{(k)}) = \underbrace{\mathbf{M}(\boldsymbol{\theta}^*)}_{=\boldsymbol{\theta}^*} + \left. \frac{\partial \mathbf{M}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} (\hat{\boldsymbol{\theta}}^{(k)} - \boldsymbol{\theta}^*) + o(\|\hat{\boldsymbol{\theta}}^{(k)} - \boldsymbol{\theta}^*\|).$$

Thus

$$\hat{\boldsymbol{\theta}}^{(k+1)} - \boldsymbol{\theta}^* = \left. \frac{\partial \mathbf{M}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} (\hat{\boldsymbol{\theta}}^{(k)} - \boldsymbol{\theta}^*) + o(\|\hat{\boldsymbol{\theta}}^{(k)} - \boldsymbol{\theta}^*\|) \quad (84)$$

and the Jacobi matrix $\left. \frac{\partial \mathbf{M}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}$ measures approximately the rate of convergence. It can be shown that

$$\left. \frac{\partial \mathbf{M}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} = [I_n^C(\boldsymbol{\theta}^*)]^{-1} I_n^{mis}(\boldsymbol{\theta}^*), \quad (85)$$

where

$$I_n^C(\boldsymbol{\theta}) = -\mathbf{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2 \ell_n^C(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \middle| \mathbb{Y}_{obs} \right]$$

can be considered as the empirical Fisher information matrix from the complete data and

$$I_n^{mis}(\boldsymbol{\theta}) = -\mathbf{E}_{\boldsymbol{\theta}} \left[\frac{\partial^2 \log f(\mathbb{Y}_{mis} | \mathbb{Y}_{obs}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \middle| \mathbb{Y}_{obs} \right],$$

can be considered as the empirical Fisher information matrix of the contribution of the missing data (that is not explained by the observed data).

Note that by (84) and (85) in the presence of missing data the convergence is only linear. Further the bigger proportion of missing data the ‘bigger’ $I_n^{mis}(\boldsymbol{\theta})$ and the slower is the convergence.

6.4 The EM algorithm in exponential families

Let the complete data \mathbb{Y} have a density with respect to a σ -finite measure μ given by

$$f(\mathbf{y}; \boldsymbol{\theta}) = \exp\left\{\sum_{j=1}^p a_j(\boldsymbol{\theta}) T_j(\mathbf{y})\right\} b(\boldsymbol{\theta}) c(\mathbf{y}) \quad (86)$$

and the standard choice of the parametric space is

$$\Theta = \left\{ \boldsymbol{\theta} : \int \exp\left\{\sum_{j=1}^p a_j(\boldsymbol{\theta}) T_j(\mathbf{y})\right\} c(\mathbf{y}) d\mu(\mathbf{y}) < \infty \right\}.$$

Note that $\mathbf{T}(\mathbb{Y}) = (T_1(\mathbb{Y}), \dots, T_p(\mathbb{Y}))^\top$ is a sufficient statistic for $\boldsymbol{\theta}$.

The log-likelihood of the complete data is now given by

$$\ell_n^C(\boldsymbol{\theta}) = \sum_{j=1}^p a_j(\boldsymbol{\theta}) T_j(\mathbb{Y}) + \log b(\boldsymbol{\theta}) + \text{const.},$$

which yields that the function Q from the EM-algorithm is given by

$$\begin{aligned} Q\left(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}\right) &= \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} \left[\ell_n^C(\boldsymbol{\theta}) | \mathbb{Y}_{obs} \right] = \sum_{j=1}^p a_j(\boldsymbol{\theta}) \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} \left[T_j(\mathbb{Y}) | \mathbb{Y}_{obs} \right] + \log b(\boldsymbol{\theta}) + \text{const.} \\ &= \sum_{j=1}^p a_j(\boldsymbol{\theta}) \widehat{T}_j^{(k)} + \log b(\boldsymbol{\theta}) + \text{const.}, \end{aligned}$$

where we put $\widehat{T}_j^{(k)} = \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} \left[T_j(\mathbb{Y}) | \mathbb{Y}_{obs} \right]$.

The nice thing about exponential families is that in the E-step of the algorithm we do not need to calculate $Q\left(\boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{(k)}\right)$ for each $\boldsymbol{\theta}$ separately but it is sufficient to calculate

$$\widehat{T}_j^{(k)} = \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} \left[T_j(\mathbb{Y}) | \mathbb{Y}_{obs} \right], \quad j = 1, \dots, p,$$

and in the M-step we maximize

$$\widehat{\boldsymbol{\theta}}^{(k+1)} = \arg \max_{\boldsymbol{\theta} \in \Theta} \left\{ \sum_{j=1}^p a_j(\boldsymbol{\theta}) \widehat{T}_j^{(k)} + \log b(\boldsymbol{\theta}) \right\}. \quad (87)$$

Interval censoring

Let $-\infty = d_0 < d_1 < \dots < d_M = \infty$ be a division of \mathbb{R} . Further let Y_1, \dots, Y_n be independent and identically distributed random variables whose exact values are not observed. Instead of each Y_i we only know that $Y_i \in (d_{q_i-1}, d_{q_i}]$, for some $q_i \in \{1, \dots, M\}$. Thus we observed independent and identically distributed random variables X_1, \dots, X_n such that $X_i = q_i$ if $Y_i \in (d_{q_i-1}, d_{q_i}]$.

Suppose now that Y_i has a density $f(y; \boldsymbol{\theta})$ of the form

$$f(y; \boldsymbol{\theta}) = \exp\left\{\sum_{j=1}^p a_j(\boldsymbol{\theta}) t_j(y)\right\} b_1(\boldsymbol{\theta}) c_1(y).$$

Thus the joint density of the random sample Y_1, \dots, Y_n is of the form (86) where

$$T_j(\mathbb{Y}) = \sum_{i=1}^n t_j(Y_i), \quad j = 1, \dots, p.$$

Thus in the E-step of the EM-algorithm it is sufficient to calculate

$$\widehat{T}_j^{(k)} = \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} [T_j(\mathbb{Y}) \mid X_1, \dots, X_n] = \sum_{i=1}^n \mathbb{E}_{\widehat{\boldsymbol{\theta}}^{(k)}} [t_j(Y_i) \mid X_i], \quad j = 1, \dots, p,$$

and the M-step is given by (87) where $b(\boldsymbol{\theta}) = b_1^n(\boldsymbol{\theta})$.

Example 47. Suppose that $Y_i \sim \text{Exp}(\lambda)$, i.e. $f(y; \lambda) = \lambda e^{-\lambda y} \mathbb{I}\{y > 0\}$. Thus $p = 1$, $t_1(y) = y$, $a_1(\lambda) = -\lambda$ and $b_1(\lambda) = \lambda$.

In the E-step one needs to calculate $\mathbb{E}_{\widehat{\lambda}^{(k)}} [Y_i \mid X_i]$. Note that the conditional distribution of Y_i given that $Y_i \in (a, b]$ has a density $\frac{\lambda e^{-\lambda y}}{e^{-\lambda a} - e^{-\lambda b}} \mathbb{I}\{y \in (a, b]\}$. Thus with the help of the integration by parts

$$\begin{aligned} \widehat{Y}_i^{(k)} &:= \mathbb{E}_{\widehat{\lambda}^{(k)}} [Y_i \mid X_i = q_i] = \frac{1}{e^{-\widehat{\lambda}^{(k)} d_{q_{i-1}}} - e^{-\widehat{\lambda}^{(k)} d_{q_i}}} \int_{d_{q_{i-1}}}^{d_{q_i}} x \widehat{\lambda}^{(k)} e^{-\widehat{\lambda}^{(k)} x} dx \\ &= \frac{d_{q_{i-1}} e^{-\widehat{\lambda}^{(k)} d_{q_{i-1}}} - d_{q_i} e^{-\widehat{\lambda}^{(k)} d_{q_i}}}{e^{-\widehat{\lambda}^{(k)} d_{q_{i-1}}} - e^{-\widehat{\lambda}^{(k)} d_{q_i}}} + \frac{1}{\widehat{\lambda}^{(k)}} \end{aligned}$$

and with the help of (87) one gets that

$$\widehat{\lambda}^{(k+1)} = \arg \max_{\lambda > 0} \left\{ Q\left(\lambda, \widehat{\lambda}^{(k)}\right) \right\} = \arg \max_{\lambda > 0} \left\{ -\lambda \sum_{i=1}^n \widehat{Y}_i^{(k)} + n \log \lambda \right\} = \frac{1}{\frac{1}{n} \sum \widehat{Y}_i^{(k)}}.$$

Literature: McLachlan and Krishnan [2008] Chapters 1.4.3, 1.5.1, 1.5.3, 2.4, 2.7, 3.2, 3.4.4, 3.5.3, 3.9 and 5.9.

7 Missing data*

For $i = 1, \dots, I$ let $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top$ represent the data of the i -th subject that could be ideally observed. Let $\mathbf{R}_i = (R_{i1}, \dots, R_{in_i})^\top$, where

$$R_{ij} = \begin{cases} 1, & \text{if } Y_{ij} \text{ is observed,} \\ 0, & \text{otherwise.} \end{cases}$$

* *Chybějící data*

Let \mathbb{Y}_{obs} represent Y_{ij} such that $R_{ij} = 1$ and \mathbb{Y}_{mis} represent Y_{ij} such that $R_{ij} = 0$. Thus the available data are given by

$$(\mathbb{Y}_{obs}, \mathbf{R}_1, \dots, \mathbf{R}_I) = (\mathbb{Y}_{obs}, \mathbf{R}),$$

where $\mathbf{R} = (\mathbf{R}_1, \dots, \mathbf{R}_I)$. Note that the complete data can be represented as

$$(\mathbf{Y}_1, \dots, \mathbf{Y}_I, \mathbf{R}) = (\mathbb{Y}_{obs}, \mathbb{Y}_{mis}, \mathbf{R}) =: (\mathbb{Y}, \mathbf{R}).$$

Suppose that the distribution of \mathbb{Y} depends on a parameter $\boldsymbol{\theta}$ (which we are interested in) and the conditional distribution of \mathbf{R} given \mathbb{Y} depends on $\boldsymbol{\psi}$. Then the joint density of the complete data can be written as

$$f(\mathbf{y}, \mathbf{r}; \boldsymbol{\theta}, \boldsymbol{\psi}) = f(\mathbf{r}|\mathbf{y}; \boldsymbol{\psi}) f(\mathbf{y}; \boldsymbol{\theta}).$$

Now integrating the above density with respect to \mathbf{y}_{mis} yields the density of the available data

$$f(\mathbf{y}_{obs}, \mathbf{r}; \boldsymbol{\theta}, \boldsymbol{\psi}) = \int f(\mathbf{y}_{obs}, \mathbf{y}_{mis}; \boldsymbol{\theta}) f(\mathbf{r}|\mathbf{y}_{obs}, \mathbf{y}_{mis}; \boldsymbol{\psi}) d\mu(\mathbf{y}_{mis}). \quad (88)$$

In what follows we will say that the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\psi}$ are *separable* if $\boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$, $\boldsymbol{\psi} \in \Omega_{\boldsymbol{\psi}}$ and $(\boldsymbol{\theta}, \boldsymbol{\psi})^T \in \Omega_{\boldsymbol{\theta}} \times \Omega_{\boldsymbol{\psi}}$.

7.1 Basic concepts for the mechanism of missing

Depending on what can be assumed about the conditional distribution of \mathbf{R} given \mathbb{Y} we distinguish three situations.

Missing completely at random (MCAR). Suppose that \mathbf{R} is independent of \mathbb{Y} , thus one can write $f(\mathbf{r}|\mathbf{y}; \boldsymbol{\psi}) = f(\mathbf{r}; \boldsymbol{\psi})$ and with the help of (88) one gets

$$f(\mathbf{y}_{obs}, \mathbf{r}; \boldsymbol{\theta}, \boldsymbol{\psi}) = f(\mathbf{y}_{obs}; \boldsymbol{\theta}) f(\mathbf{r}; \boldsymbol{\psi}),$$

which further implies that the observed log-likelihood is of the form

$$\ell_{obs}(\boldsymbol{\theta}, \boldsymbol{\psi}) = \log f(\mathbb{Y}_{obs}; \boldsymbol{\theta}) + \log f(\mathbf{R}; \boldsymbol{\psi}).$$

Note that if the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\psi}$ are separable then the second term on the right-hand side of the above equation does not depend on $\boldsymbol{\theta}$ and can be ignored when one is interested only in $\boldsymbol{\theta}$.

Example 48. Let Y_1, \dots, Y_n be a random sample from the exponential distribution $\text{Exp}(\lambda)$. Let R_1, \dots, R_n be a random sample independent with Y_1, \dots, Y_n and R_i follows a Bernoulli distribution with a parameter p_i (e.g. $p_i = \frac{1}{1+i}$).

Missing at random (MAR). Suppose that the conditional distribution of \mathbf{R} given \mathbb{Y} is the same as the conditional distribution of \mathbf{R} given \mathbb{Y}_{obs} . Thus one can write $f(\mathbf{r}|\mathbf{y}; \boldsymbol{\psi}) = f(\mathbf{r}|\mathbf{y}_{obs}; \boldsymbol{\psi})$ and with the help of (88)

$$f(\mathbf{y}_{obs}, \mathbf{r}; \boldsymbol{\theta}, \boldsymbol{\psi}) = f(\mathbf{y}_{obs}; \boldsymbol{\theta})f(\mathbf{r}|\mathbf{y}_{obs}; \boldsymbol{\psi}),$$

which further implies that the observed log-likelihood is of the form

$$\ell_{obs}(\boldsymbol{\theta}, \boldsymbol{\psi}) = \log f(\mathbb{Y}_{obs}; \boldsymbol{\theta}) + \log f(\mathbf{R}|\mathbb{Y}_{obs}; \boldsymbol{\psi}).$$

Note that although MAR is not so strict in assumptions as MCAR, also here the second term on the right-hand side of the above equation does not depend on $\boldsymbol{\theta}$ provided that $\boldsymbol{\theta}$ and $\boldsymbol{\psi}$ are separable.

Example 49. Let $(\mathbf{X}_1^\top, Y_1, R_1)^\top, \dots, (\mathbf{X}_n^\top, Y_n, R_n)^\top$ be independent and identically distributed random vectors, where the covariates $\mathbf{X}_1, \dots, \mathbf{X}_n$ are always completely observed. Let R_i stand for the indicator of missing of Y_i and

$$\mathbb{P}(R_i = 1 | \mathbf{X}_i, Y_i) = r(\mathbf{X}_i),$$

where $r(\mathbf{x})$ is a given (but possibly unknown) function.

Missing not at random (MNAR). In this concept neither the distribution of \mathbf{R} is independent of \mathbb{Y} nor the conditional distribution of \mathbf{R} given \mathbb{Y}_{obs} is independent of \mathbb{Y}_{mis} . Thus the density of the observed data is generally given by (88). To proceed one has to make some other assumptions about the conditional distribution of \mathbf{R} given \mathbb{Y} (i.e. about the density $f(\mathbf{r}|\mathbf{y}_{obs}, \mathbf{y}_{mis}; \boldsymbol{\psi})$).

Example 50. *Maximum likelihood estimator for the right-censored data from an exponential distribution.* Suppose that Y_1, \dots, Y_n is a random sample from the exponential distribution with the density $f(x; \lambda) = \lambda e^{-\lambda x} \mathbb{I}\{x > 0\}$. Nevertheless we observe Y_i only if $Y_i \leq C$, where C is a known constant (e.g. duration of the study). If $Y_i > C$ then we do not observe the value of Y_i (we only know that Y_i is greater than C).

Note that

$$f(\mathbf{y}_{obs}, \mathbf{y}_{mis}; \lambda) = \prod_{i=1}^n \lambda e^{-\lambda y_i}$$

and

$$f(\mathbf{r}|\mathbf{y}_{obs}, \mathbf{y}_{mis}) = \prod_{i=1}^n [\mathbb{I}\{y_i \leq C\}]^{r_i} [\mathbb{I}\{y_i > C\}]^{1-r_i}.$$

Although this conditional density depends on \mathbf{y}_{mis} (thus we are in a situation of MNAR), we can proceed because this conditional density is completely known.

Let n_0 be the number of fully observed Y_i (i.e. $n_0 = \sum_{i=1}^n \mathbb{1}\{Y_i \leq C\}$). For simplicity of notation assume that Y_1, \dots, Y_n are ordered in such a way that Y_1, \dots, Y_{n_0} are fully observed and Y_{n_0+1}, \dots, Y_n are censored (i.e. $Y_i > C$ for $i \in \{n_0 + 1, \dots, n\}$). Thus the corresponding components of \mathbf{R} are given by $R_i = 1$ for $i \in \{1, \dots, n_0\}$ and zero otherwise.

Now with the help of (88) one can calculate

$$\begin{aligned} f(\mathbf{Y}_{obs}, \mathbf{R}; \lambda) &= \prod_{i=1}^{n_0} \lambda e^{-\lambda Y_i} \int_C^\infty \cdots \int_C^\infty \prod_{i=n_0+1}^n \lambda e^{-\lambda y_i} dy_{n_0+1}, \dots, dy_n \\ &= \lambda^{n_0} e^{-\lambda \sum_{i=1}^{n_0} Y_i} [e^{-\lambda C}]^{n-n_0}. \end{aligned}$$

The corresponding log-likelihood of the observed data is

$$\ell_{obs}(\lambda) = n_0 \log \lambda - \lambda \sum_{i=1}^{n_0} Y_i - (n - n_0)C\lambda,$$

which is maximised at

$$\hat{\lambda}_n = \frac{1}{\frac{1}{n_0} \sum_{i=1}^{n_0} Y_i + \frac{(n-n_0)C}{n_0}}.$$

Note that the above example is in fact rather exceptional as the missing mechanism is given by the design of the study and thus known.

The general problem of all the concepts is that if missing is not a part of the design of the study then no assumptions about the relationship of Y_{mis} and \mathbf{R} can be verified as we do not observe Y_{mis} .

The expected
end of class 23
(10.5.2024)

7.2 Methods for dealing with missing data

Complete case analysis (CCA)

In the analysis we use only the subjects with the full record, i.e. only subjects for which no information is missing.

Advantages and disadvantages:

- + simplicity;
- the inference about θ is ‘biased’ (i.e. the parameter θ is generally not identified), if MCAR does not hold;
- even if MCAR holds, then this method may not provide an effective use of data.

Example 51. Suppose that we have five observations on each subject. Each observation is missing with probability 0.1 and the observations are missing independently on each other. Thus on average only 59% ($0.9^5 \doteq 0.59$) of the records will be complete.

Available case analysis (ACA)

In each of the analyses one uses all the data that are available for this particular analysis.

Example 52. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a random sample from $N((\mu_1, \mu_2, \mu_3)^\top, \Sigma_{3 \times 3})$. Then the covariance $\sigma_{ij} = \text{cov}(X_{1i}, X_{1j})$ is estimated from all the vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ for which both the i -th and the j -th coordinate is observed.

Advantages and disadvantages:

- + simplicity;
- + more data can be used than with CCA;
- the inference about $\boldsymbol{\theta}$ is biased, if MCAR does not hold;
- it can result in estimates with strange features (e.g. there is no guarantee that the estimate of the variance matrix $\hat{\Sigma}$ in Example 52 is positive semidefinite).

Direct (ignorable) observed likelihood

The inference is based on $\log f(\mathcal{Y}_{obs}; \boldsymbol{\theta})$, that is the distribution of \mathbf{R} is ‘ignored’.

Advantages and disadvantages:

- + If the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\psi}$ are separable then this method is not biased and does not lose any information provided MAR holds;
- The observed log-likelihood $\ell_{obs}(\boldsymbol{\theta})$ might be difficult to calculate. Nevertheless, sometimes the EM algorithm can be helpful.

Imputation

In this method the missing observations are estimated (‘imputed’) and then one works with the data as if there were no missing values.

Advantages and disadvantages:

- + If the missing values are estimated appropriately, it can give ‘reasonable’ estimates of the unknown parameters;
- + One can use the *completed* dataset also for other analyses;
- The standard estimates of the (asymptotic) variances of the estimates of the parameters computed from the completed dataset are too optimistic (too low). The reason is that an appropriate estimate of variance should reflect that part of the data has been imputed.

Example 53. Suppose that X_1, \dots, X_n is a random sample. Further suppose that we observe only X_1, \dots, X_{n_0} for some $n_0 < n$ and the remaining observations X_{n_0+1}, \dots, X_n are missing. For $i = n_0 + 1, \dots, n$ let the missing observations be estimated as $\hat{X}_i = \frac{1}{n_0} \sum_{j=1}^{n_0} X_j$. Then the standard estimate of $\mu = \mathbb{E} X_1$ is given by

$$\hat{\mu}_n = \frac{1}{n} \left(\sum_{i=1}^{n_0} X_i + \sum_{i=n_0+1}^n \hat{X}_i \right) = \frac{1}{n_0} \sum_{j=1}^{n_0} X_j$$

and seems to be reasonable.

But the standard estimate of the variance of $\hat{\mu}_n$ computed from the completed dataset

$$\widehat{\text{var}}(\hat{\mu}_n) = \frac{S_n^2}{n}, \quad \text{where} \quad S_n^2 = \frac{1}{n-1} \left(\sum_{i=1}^{n_0} (X_i - \hat{\mu}_n)^2 + \sum_{i=n_0+1}^n (\hat{X}_i - \hat{\mu}_n)^2 \right)$$

is too small. The first reason is that S_n^2 as the estimate of $\text{var}(X_1)$ is

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^{n_0} (X_i - \hat{\mu}_n)^2 = \frac{n_0-1}{n-1} S_{n_0}^2 < S_{n_0}^2.$$

The second reason is that the factor $\frac{1}{n}$ assumes that there are n independent observations, but in fact there are only n_0 independent observations.

The expected
end of class 24
(17.5.2024)

Multiple imputation

In this method the missing observations are imputed several times. Formally, for $j = 1, \dots, M$ let $\hat{Y}_{mis}^{(j)}$ be the imputed values in the j -th round. Further let $\hat{\theta}_j$ be the estimate of the parameter θ from the completed data $(Y_{obs}, \hat{Y}_{mis}^{(j)})$. Then the final estimate of the parameter θ is given by

$$\hat{\theta}_{MI} = \frac{1}{M} \sum_{j=1}^M \hat{\theta}_j.$$

The advantage of this method is that one can also estimate the (asymptotic) variance of this estimator by

$$\widehat{\text{var}}(\hat{\theta}_{MI}) = \bar{V}_M + \left(1 + \frac{1}{M}\right) \mathbb{B}_M, \quad (89)$$

where

$$\bar{V}_M = \frac{1}{M} \sum_{j=1}^M \hat{V}_j \quad \text{and} \quad \mathbb{B}_M = \frac{1}{M-1} \sum_{j=1}^M \left(\hat{\theta}_j - \hat{\theta}_{MI} \right) \left(\hat{\theta}_j - \hat{\theta}_{MI} \right)^\top,$$

with \hat{V}_j being a standard estimate of the asymptotic variance calculated from the completed data $\hat{Y}^{(j)} = (Y_{obs}, \hat{Y}_{mis}^{(j)})$.

The rationale of the formula (89) is as follows. Note that

$$\text{var}(\widehat{\boldsymbol{\theta}}_{MI}) = \text{E}(\text{var}(\widehat{\boldsymbol{\theta}}_{MI} | \widehat{\mathbf{Y}}^{(j)})) + \text{var}(\text{E}(\widehat{\boldsymbol{\theta}}_{MI} | \widehat{\mathbf{Y}}^{(j)})).$$

Now the first term on right-hand side of the above equation is estimated by $\overline{\mathbb{V}}_M$ and the second term is estimated by \mathbb{B}_M .

Example 54. In Example 53 one can for instance impute the values X_{n_0+1}, \dots, X_n by a random sample from $\text{N}(\widehat{\mu}, \widehat{\sigma}^2)$, where $\widehat{\mu} = \overline{X}_{n_0}$ and $\widehat{\sigma}^2 = S_{n_0}^2$ are the sample mean and variance calculated from the observed data. Put $\widehat{V}_j = \frac{S_n^{2(j)}}{n}$, where $S_n^{2(j)}$ is the sample variance calculated from the j -th completed sample. Then one can show that

$$\overline{\mathbb{V}}_M \xrightarrow[M \rightarrow \infty]{\text{a.s.}} \frac{S_{n_0}^2}{n}. \quad (90)$$

Further let $\widehat{\theta}_j = \overline{Y}_n^{(j)}$ be the sample mean calculated from the j -th completed sample. Then it can be shown that

$$B_M \xrightarrow[M \rightarrow \infty]{\text{a.s.}} \frac{S_{n_0}^2(n - n_0)}{n^2}. \quad (91)$$

Now combining (90) and (91) yields that

$$\overline{\mathbb{V}}_M + B_M \xrightarrow[M \rightarrow \infty]{\text{a.s.}} S_{n_0}^2 \left(\frac{2}{n} - \frac{n_0}{n^2} \right).$$

Further it is straightforward to prove that for $n_0 < n$

$$S_{n_0}^2 \left(\frac{2}{n} - \frac{n_0}{n^2} \right) < \frac{S_{n_0}^2}{n_0},$$

where the right-hand side of the above inequality represents the standard estimate of the variance of \overline{X}_{n_0} (that assumes MCAR). This indicates that when doing multiple imputation, one needs to take into consideration also the variability that comes from the fact that one uses the estimates $\widehat{\mu}, \widehat{\sigma}^2$ instead of the true values of μ and σ . This can be done very naturally within the framework of Bayesian statistics.

Advantages and disadvantages:

- + If the missing values are estimated appropriately, it can give ‘reasonable’ estimate of the unknown parameter as well as of the variance of this estimate.
- To be done properly it requires the knowledge of Bayesian approach to statistics.

Re-weighting

Roughly speaking in this method each observation is given a weight (w_i) that is proportional to the inverse probability of being observed (π_i), i.e.

$$w_i = \frac{\frac{1}{\pi_i}}{\sum_{j:R_j=1} \frac{1}{\pi_j}}, \quad i \in \{j : R_j = 1\}.$$

All the procedures are now weighted with respect to these weights, e.g. the M -estimator of a parameter θ is given by

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \sum_{i:R_i=1} w_i \rho(\mathbf{X}_i; \theta).$$

Example 55. Suppose we have a study where for a large number of patients some basic and cheap measurements have been done resulting in $\mathbf{Z}_1, \dots, \mathbf{Z}_N$. Now a random subsample \mathbb{S} of size n from these patients has been done for some more expensive measurements. Note that then $\mathbb{S} = \{j \in \{1, \dots, N\} : R_j = 1\}$.

This method can be also used where one has some auxiliary variables $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ that can be used to estimate the probabilities π_i with the help of for instance a logistic regression.

The expected
end of class 25
(21. 5. 2024)

Appendix

Inverse function theorem

The following theorem is sometimes also called the theorem about the local diffeomorphism. It follows easily from the implicit function theorem applied to the function $\mathbf{g}(\mathbf{x}, \mathbf{y}) = \mathbf{x} - \mathbf{f}(\mathbf{y})$.

Theorem A13. *Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ have continuous first order partial derivatives in a neighbourhood of the point $\mathbf{a} \in \mathbb{R}^n$ and the Jacobi matrix $\mathbb{D}_{\mathbf{f}}(\mathbf{a})$ is a non-singular matrix. Then there exist open neighbourhoods U of the point \mathbf{a} and V of the point $\mathbf{f}(\mathbf{a})$ such that \mathbf{f} is a bijection of U on V . Further there exists an inverse function \mathbf{f}^{-1} on V with the continuous first order partial derivatives.*

Lemma about the distribution of a quadratic form

The following lemma can be found as Theorem 4.16 in [Anděl \[2007\]](#).

Lemma A6. *Let $\mathbf{Z} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbb{V})$, where \mathbb{V} is $p \times p$ matrix. Let \mathbb{B} be a positively semidefinite matrix such that $\mathbb{B}\mathbb{V}$ is an idempotent (nonzero) matrix. Then $\mathbf{Z}^T \mathbb{B} \mathbf{Z} \sim \chi_{\text{tr}(\mathbb{B}\mathbb{V})}^2$.*

Banach fixed point theorem

Definition. Let (P, ρ) be a metric space. Then a map $T : P \mapsto P$ is called a *contraction mapping* on P if there exists $q \in [0, 1)$ such that for all $x, y \in P$

$$\rho(T(x), T(y)) \leq q\rho(x, y).$$

Theorem A14. *Let (P, ρ) be a non-empty complete metric space with a contraction mapping $T : P \mapsto P$. Then T admits a unique fixed-point $x^* \in P$ (i.e. $T(x^*) = x^*$).*

Bayes theorem for densities

Theorem A15. *Suppose that $\mathbf{X} = (X_1, \dots, X_k)^T$ and $\mathbf{Z} = (Z_1, \dots, Z_G)^T$ be random vectors defined on the same probability space. Let $f_{\mathbf{X}}$ and $f_{\mathbf{Z}}$ be the densities of \mathbf{X} and \mathbf{Z} respectively and $f_{\mathbf{X}|\mathbf{Z}}$ be the conditional density of \mathbf{X} given \mathbf{Z} . Then the conditional density of \mathbf{Z} given \mathbf{X} equals*

$$f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x}) = \begin{cases} \frac{f_{\mathbf{X}|\mathbf{Z}}(\mathbf{x}|\mathbf{z})f_{\mathbf{Z}}(\mathbf{z})}{f_{\mathbf{X}}(\mathbf{x})}, & \text{for } f_{\mathbf{X}}(\mathbf{x}) > 0, \\ 0, & \text{for } f_{\mathbf{X}}(\mathbf{x}) = 0. \end{cases}$$

Proof. The proof follows from the fact that $f_{\mathbf{X}, \mathbf{Z}}(\mathbf{x}, \mathbf{z}) = f_{\mathbf{X}|\mathbf{Z}}(\mathbf{x}|\mathbf{z})f_{\mathbf{Z}}(\mathbf{z})$ is the joint density of $\begin{pmatrix} \mathbf{X} \\ \mathbf{Z} \end{pmatrix}$ and then by the definition of the conditional density. For details see e.g. Chapter 3.5 of [Anděl \[2007\]](#). □

References

- J. Anděl. *Základy matematické statistiky*. Matfyzpress, Praha, 2007.
- N. L. Hjort and D. Pollard. Asymptotics for minimisers of convex processes. *arXiv preprint, arXiv:1107.3806*, 2011.
- B. Hoadley. Asymptotic properties of maximum likelihood estimators for the independent not identically distributed case. *Annals of Mathematical Statistics*, 42:1977–1991, 1971.
- J. Jiang. *Large sample techniques for statistics*. Springer Texts in Statistics. Springer, New York, 2010. ISBN 978-1-4419-6826-5. doi: 10.1007/978-1-4419-6827-2. URL <http://dx.doi.org/10.1007/978-1-4419-6827-2>.
- R. Koenker. *Quantile Regression*. Cambridge University Press, New York, 2005.
- M. Kulich. Maximum likelihood estimation theory, 2014. URL <https://www.karlin.mff.cuni.cz/~kulich/vyuka/glm/index.html>.
- P. Lachout. *Teorie pravděpodobnosti*. Karolinum, 2004. Skripta.
- H. Leeb and B. M. Pötscher. Sparse estimators and the oracle property, or the return of hedges' estimator. *Journal of Econometrics*, 142(1):201–211, 2008.
- E. L. Lehmann and G. Casella. *Theory of point estimation*. Springer, New York, 1998.
- R. A. Maronna, R. D. Martin, and V. J. Yohai. *Robust statistics*. Wiley, Chichester, 2006.
- G. McLachlan and T. Krishnan. *The EM algorithm and extensions*. Wiley, New York, 2008. Second Edition.
- S. Nagy. Mathematical statistics 2. Course notes for NMSA332. URL <https://www.karlin.mff.cuni.cz/~nagy/NMSA332/NMSA332.pdf>.
- Y. Pawitan. *In all likelihood: statistical modelling and inference using likelihood*. Oxford University Press, 2001.
- P. K. Sen, J. M. Singer, and A. C. P. de Lima. *From finite sample to asymptotic methods in statistics*. Cambridge University Press, 2010.
- R. J. Serfling. *Approximation Theorems of Mathematical Statistics*. Wiley, New York, 1980.
- A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, New York, 2000. ISBN 0521784506.

- H. White. A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica*, 48:817–838, 1980.
- H. White. Maximum likelihood estimation of misspecified models. *Econometrica*, 50:1–25, 1982.
- C. J. Wu. On the convergence properties of the EM algorithm. *Annals of Statistics*, 11: 95–103, 1983.
- K. Zvára. *Regrese*. MATFYZPRESS, 2008.