

Bias Correction in Maximum Likelihood Logistic Regression (Schaefer, R. 1983)

Jens Rahnfeld

SS 2021

- Maximum Likelihood Schätzer ist nur asymptotisch unverzerrt.

- Maximum Likelihood Schätzer ist nur asymptotisch unverzerrt.
- Bei kleineren Datensätzen kann es zu einem hohen Bias kommen.

- Maximum Likelihood Schätzer ist nur asymptotisch unverzerrt.
- Bei kleineren Datensätzen kann es zu einem hohen Bias kommen.
- Idee: Korrigiere den Maximum Likelihood Schätzer um dessen Bias

- Datensatz $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ aus n i.i.d. Datenpunkten (x_i, y_i) , mit zugrundeliegender Verteilung:

$$\mathbb{P}(y = 1 \mid x') = \frac{\exp(\beta_t^\top x)}{1 + \exp(\beta_t^\top x)}$$

wobei $x = \begin{pmatrix} 1 \\ x' \end{pmatrix} \in \mathbb{R}^p$, $y \in \{0, 1\}$, $\beta_t \in \mathbb{R}^p$.

Logistisches Regressionsmodell

- Datensatz $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ aus n i.i.d. Datenpunkten (x_i, y_i) , mit zugrundeliegender Verteilung:

$$\mathbb{P}(y = 1 \mid x') = \frac{\exp(\beta_t^\top x)}{1 + \exp(\beta_t^\top x)}$$

wobei $x = \begin{pmatrix} 1 \\ x' \end{pmatrix} \in \mathbb{R}^p$, $y \in \{0, 1\}$, $\beta_t \in \mathbb{R}^p$.

- β_t bezeichnet wahren Parameter und ist konstant.

- Datensatz $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ aus n i.i.d. Datenpunkten (x_i, y_i) , mit zugrundeliegender Verteilung:

$$\mathbb{P}(y = 1 \mid x') = \frac{\exp(\beta_t^\top x)}{1 + \exp(\beta_t^\top x)}$$

wobei $x = \begin{pmatrix} 1 \\ x' \end{pmatrix} \in \mathbb{R}^p$, $y \in \{0, 1\}$, $\beta_t \in \mathbb{R}^p$.

- β_t bezeichnet wahren Parameter und ist konstant.
- Wir schreiben x_i für $\begin{pmatrix} 1 \\ x_i \end{pmatrix} = \left(1 \quad x_i^{(2)} \quad \dots \quad x_i^{(p)}\right)^\top$

- Datensatz:

$$\mathbf{X} = \begin{pmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(2)} & \dots & x_1^{(p)} \\ \vdots & & & \\ 1 & x_n^{(2)} & \dots & x_n^{(p)} \end{pmatrix} \in \mathbb{R}^{n \times p}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

- Datensatz:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(2)} & \dots & x_1^{(p)} \\ \vdots & & & \\ 1 & x_n^{(2)} & \dots & x_n^{(p)} \end{pmatrix} \in \mathbb{R}^{n \times p}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

- Wahrscheinlichkeiten:

$$p_i(\beta) = \mathbb{P}_\beta(y_i = 1 \mid x_i) = \frac{\exp(\beta^\top x_i)}{1 + \exp(\beta^\top x_i)}, \quad \mathbf{p}(\beta) = \begin{pmatrix} p_1(\beta) \\ \vdots \\ p_n(\beta) \end{pmatrix}$$

- Datensatz:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(2)} & \dots & x_1^{(p)} \\ \vdots & & & \\ 1 & x_n^{(2)} & \dots & x_n^{(p)} \end{pmatrix} \in \mathbb{R}^{n \times p}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

- Wahrscheinlichkeiten:

$$p_i(\beta) = \mathbb{P}_\beta(y_i = 1 \mid x_i) = \frac{\exp(\beta^\top x_i)}{1 + \exp(\beta^\top x_i)}, \quad \mathbf{p}(\beta) = \begin{pmatrix} p_1(\beta) \\ \vdots \\ p_n(\beta) \end{pmatrix}$$

- wahre Varianz:

$$\mathbf{V} = \text{diag}(v_1, \dots, v_n) = \begin{pmatrix} p_1(\beta_t)(1 - p_1(\beta_t)) & & & \\ & \ddots & & \\ & & \ddots & \\ & & & p_n(\beta_t)(1 - p_n(\beta_t)) \end{pmatrix}$$

We assume that the independent variables are bounded so that $\lim (\mathbf{X}^T \mathbf{V} \mathbf{X})/n$ is finite and positive definite as $n \rightarrow \infty$. Then, under certain regularity conditions, it is well known³ that $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ converges in distribution to a p -dimensional multivariate normal distribution with mean vector $\mathbf{0}$ and covariance matrix $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. Thus if the sample size is large, $\hat{\boldsymbol{\beta}}$ is approximately unbiased and has approximate covariance matrix given by $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. If the sample size is small, however, $\hat{\boldsymbol{\beta}}$ might be biased.

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- \mathbf{X} wird im Paper als deterministische Matrix betrachtet.

We assume that the independent variables are bounded so that $\lim (\mathbf{X}^T \mathbf{V} \mathbf{X})/n$ is finite and positive definite as $n \rightarrow \infty$. Then, under certain regularity conditions, it is well known³ that $(\hat{\beta} - \beta)$ converges in distribution to a p -dimensional multivariate normal distribution with mean vector $\mathbf{0}$ and covariance matrix $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. Thus if the sample size is large, $\hat{\beta}$ is approximately unbiased and has approximate covariance matrix given by $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. If the sample size is small, however, $\hat{\beta}$ might be biased.

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- \mathbf{X} wird im Paper als deterministische Matrix betrachtet.
- Beachte:

$$\mathbf{X}^T \mathbf{V} \mathbf{X} = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} v_1 \cdot x_1^T \\ \vdots \\ v_n \cdot x_n^T \end{pmatrix} = \sum_{i=1}^n v_i \cdot x_i x_i^T$$

We assume that the independent variables are bounded so that $\lim (\mathbf{X}^T \mathbf{V} \mathbf{X})/n$ is finite and positive definite as $n \rightarrow \infty$. Then, under certain regularity conditions, it is well known³ that $(\hat{\beta} - \beta)$ converges in distribution to a p -dimensional multivariate normal distribution with mean vector $\mathbf{0}$ and covariance matrix $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. Thus if the sample size is large, $\hat{\beta}$ is approximately unbiased and has approximate covariance matrix given by $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$. If the sample size is small, however, $\hat{\beta}$ might be biased.

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- \mathbf{X} wird im Paper als deterministische Matrix betrachtet.
- Beachte:

$$\mathbf{X}^T \mathbf{V} \mathbf{X} = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} v_1 \cdot x_1^T \\ \vdots \\ v_n \cdot x_n^T \end{pmatrix} = \sum_{i=1}^n v_i \cdot x_i x_i^T$$

- Als Zufallsvariable betrachtet, fordern wir, dass $E[v \cdot x x^T]$ endlich, positiv definit und $\mathbf{X}^T \mathbf{V} \mathbf{X}$ fast sicher invertierbar definit ist. Weiter fordern wir $\|x\| \leq C$ fast sicher für ein $C < \infty$.

- Maximiere Log-Likelihood:

$$l(\beta) = \sum_{i=1}^n y_i \cdot \beta^\top x_i - \log(1 + \exp(\beta^\top x_i))$$

- Maximiere Log-Likelihood:

$$l(\beta) = \sum_{i=1}^n y_i \cdot \beta^\top x_i - \log(1 + \exp(\beta^\top x_i))$$

- Maximierer $\hat{\beta}$ ist Lösung von

$$\dot{l}(\beta) = \mathbf{0}$$

wobei

$$\dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta)) \cdot x_i = \mathbf{X}^\top (\mathbf{y} - \mathbf{p}(\beta))$$

$$\dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta)) \cdot x_i$$

- Entwickle p_i im Punkt β_t :

$$\begin{aligned} \dot{p}_i(\beta) &= \frac{d}{d\beta} \frac{\exp(\beta^\top x_i)}{1 + \exp(\beta^\top x_i)} \\ &= \frac{\exp(\beta^\top x_i)(1 + \exp(\beta^\top x_i)) - \exp(\beta^\top x_i)^2}{(1 + \exp(\beta^\top x_i))^2} \cdot x_i^\top \\ &= \frac{\exp(\beta^\top x_i)}{(1 + \exp(\beta^\top x_i))^2} \cdot x_i^\top \\ &= p_i(\beta) \cdot (1 - p_i(\beta)) \cdot x_i^\top \end{aligned}$$

$$\dot{p}_i(\beta) = \underbrace{p_i(\beta) \cdot (1 - p_i(\beta))}_{=: v_i(\beta)} \cdot x_i^\top$$

- Entwickle p_i im Punkt β_t :

$$\begin{aligned} \frac{d}{d\beta} p_i(\beta)(1 - p_i(\beta)) &= v_i(\beta)(1 - p_i(\beta)) \cdot x_i^\top - p_i(\beta)v_i(\beta) \cdot x_i^\top \\ &= v_i(\beta)(1 - 2p_i(\beta)) \cdot x_i^\top \end{aligned}$$

$$\dot{p}_i(\beta) = \underbrace{p_i(\beta) \cdot (1 - p_i(\beta))}_{=: v_i(\beta)} \cdot x_i^\top$$

- Entwickle p_i im Punkt β_t :

$$\begin{aligned} \frac{d}{d\beta} p_i(\beta)(1 - p_i(\beta)) &= v_i(\beta)(1 - p_i(\beta)) \cdot x_i^\top - p_i(\beta)v_i(\beta) \cdot x_i^\top \\ &= v_i(\beta)(1 - 2p_i(\beta)) \cdot x_i^\top \end{aligned}$$

$$\begin{aligned} \Rightarrow \ddot{p}_i(\beta) &= \frac{d}{d\beta} p_i(\beta)(1 - p_i(\beta)) \cdot x_i \\ &= v_i(\beta)(1 - 2p_i(\beta)) \cdot x_i x_i^\top \end{aligned}$$

$$\dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta)) \cdot x_i$$

- Entwickle p_i im Punkt β_t :

$$\dot{p}_i(\beta) = v_i(\beta) \cdot x_i^\top$$

$$\ddot{p}_i(\beta) = v_i(\beta)(1 - 2p_i(\beta)) \cdot x_i x_i^\top$$

$$\begin{aligned} \dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta_t) - \dot{p}_i(\beta_t)(\beta - \beta_t) - \frac{1}{2} \cdot (\beta - \beta_t)^\top \ddot{p}_i(\beta_t)(\beta - \beta_t)) \cdot x_i \\ - o(\|\beta - \beta_t\|^2) x_i \end{aligned}$$

$$\dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta)) \cdot x_i$$

- Entwickle p_i im Punkt β_t :

$$\dot{p}_i(\beta) = v_i(\beta) \cdot x_i^\top$$

$$\ddot{p}_i(\beta) = v_i(\beta)(1 - 2p_i(\beta)) \cdot x_i x_i^\top$$

$$\begin{aligned} \dot{l}(\beta) = \sum_{i=1}^n (y_i - p_i(\beta_t) - \dot{p}_i(\beta_t)(\beta - \beta_t) - \frac{1}{2} \cdot (\beta - \beta_t)^\top \ddot{p}_i(\beta_t)(\beta - \beta_t)) \cdot x_i \\ - o(\|\beta - \beta_t\|^2) x_i \end{aligned}$$

- Wir vernachlässigen den $o(\|\beta - \beta_t\|^2)$ Term vorerst.

- Terme 0. Ordnung:

$$\sum_{i=1}^n (y_i - p_i(\beta_t)) \cdot x_i = \mathbf{X}^\top (\mathbf{y} - \mathbf{p}(\beta_t))$$

- Terme 1. Ordnung:

$$\sum_{i=1}^n \dot{p}_i(\beta_t) (\beta - \beta_t) x_i = \sum_{i=1}^n v_i \cdot x_i^\top (\beta - \beta_t) x_i = \mathbf{X}^\top \mathbf{V} \mathbf{X} (\beta - \beta_t)$$

- Terme 2. Ordnung:

$$\begin{aligned}
 & -\frac{1}{2} \sum_{i=1}^n (\beta - \beta_t)^\top \ddot{\mathbf{p}}_i(\beta_t) (\beta - \beta_t) \cdot \mathbf{x}_i \\
 & = \mathbf{X}^\top \cdot \left(-\frac{1}{2}\right) \begin{pmatrix} (\beta - \beta_t)^\top \ddot{\mathbf{p}}_1(\beta - \beta_t) \\ \vdots \\ (\beta - \beta_t)^\top \ddot{\mathbf{p}}_n(\beta - \beta_t) \end{pmatrix}
 \end{aligned}$$

wobei $\ddot{\mathbf{p}}_i = \ddot{p}_i(\beta_t) = v_i(1 - 2p_i) \cdot \mathbf{x}_i \mathbf{x}_i^\top \in \mathbb{R}^{p \times p}$

- Insgesamt erhalten wir:

$$\begin{aligned}
 \dot{l}(\beta) &= \sum_{i=1}^n (y_i - p_i(\beta_t) - \dot{p}_i(\beta_t)(\beta - \beta_t) - \frac{1}{2} \cdot (\beta - \beta_t)^\top \ddot{p}_i(\beta_t)(\beta - \beta_t)) \cdot x_i \\
 &= \mathbf{X}^\top \left(\mathbf{y} - \mathbf{p}(\beta_t) - \mathbf{V}\mathbf{X}(\beta - \beta_t) - \frac{1}{2} \cdot \begin{pmatrix} (\beta - \beta_t)^\top \ddot{\mathbf{p}}_1(\beta - \beta_t) \\ \vdots \\ (\beta - \beta_t)^\top \ddot{\mathbf{p}}_n(\beta - \beta_t) \end{pmatrix} \right)
 \end{aligned}$$

Logistische Regression als Latent Variable Modell

(Alternatively, one could use iterative weighted least squares on the model

$$\mathbf{y} = \boldsymbol{\pi} + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\pi}$ is the $(n \times 1)$ vector of the π_i and $\boldsymbol{\varepsilon}$ is an $(n \times 1)$ vector of independent Bernoulli random errors, ε_i , with $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Var}(\boldsymbol{\varepsilon}) = \mathbf{V} = \text{diag}\{v_i\} = \text{diag}\{\pi_i(1 - \pi_i)\}$ to arrive at the same estimate.⁵⁾

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- Relevant ist, dass $\mathbf{y} - \boldsymbol{p}(\boldsymbol{\beta}_t)$ durch $\boldsymbol{\varepsilon}$ ersetzt werden kann.

Logistische Regression als Latent Variable Modell

(Alternatively, one could use iterative weighted least squares on the model

$$\mathbf{y} = \boldsymbol{\pi} + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\pi}$ is the $(n \times 1)$ vector of the π_i and $\boldsymbol{\varepsilon}$ is an $(n \times 1)$ vector of independent Bernoulli random errors, ε_i , with $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Var}(\boldsymbol{\varepsilon}) = \mathbf{V} = \text{diag}\{v_i\} = \text{diag}\{\pi_i(1 - \pi_i)\}$ to arrive at the same estimate.⁵⁾

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- Relevant ist, dass $\mathbf{y} - \boldsymbol{p}(\boldsymbol{\beta}_t)$ durch $\boldsymbol{\varepsilon}$ ersetzt werden kann.
- Bemerkung:

$$\begin{aligned} E[y_i] &= E[E[1_{\{y_i=1\}} \mid x_i]] = E[\mathbb{P}(y_i = 1 \mid x_i)] = E[p_i(\boldsymbol{\beta}_t)] \\ E[(y_i - p_i(\boldsymbol{\beta}_t))^2] &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 \cdot 1_{\{y_i=1\}} + p_i(\boldsymbol{\beta}_t)^2 \cdot 1_{\{y_i=0\}}] \\ &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 E[1_{\{y_i=1\}} \mid x_i] + p_i(\boldsymbol{\beta}_t)^2 E[1_{\{y_i=0\}} \mid x_i]] \\ &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 p_i(\boldsymbol{\beta}_t) + p_i(\boldsymbol{\beta}_t)^2 (1 - p_i(\boldsymbol{\beta}_t))] \\ &= E[p_i(\boldsymbol{\beta}_t)(1 - p_i(\boldsymbol{\beta}_t))] \end{aligned}$$

Logistische Regression als Latent Variable Modell

(Alternatively, one could use iterative weighted least squares on the model

$$\mathbf{y} = \boldsymbol{\pi} + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\pi}$ is the $(n \times 1)$ vector of the π_i and $\boldsymbol{\varepsilon}$ is an $(n \times 1)$ vector of independent Bernoulli random errors, ε_i , with $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Var}(\boldsymbol{\varepsilon}) = \mathbf{V} = \text{diag}\{v_i\} = \text{diag}\{\pi_i(1 - \pi_i)\}$ to arrive at the same estimate.⁵⁾

Figure: Schaefer R., Bias Correction in Maximum Likelihood Logistic Regression, 1983

- Relevant ist, dass $\mathbf{y} - \mathbf{p}(\boldsymbol{\beta}_t)$ durch $\boldsymbol{\varepsilon}$ ersetzt werden kann.
- Bemerkung:

$$\begin{aligned} E[y_i] &= E[E[1_{\{y_i=1\}} \mid x_i]] = E[\mathbb{P}(y_i = 1 \mid x_i)] = E[p_i(\boldsymbol{\beta}_t)] \\ E[(y_i - p_i(\boldsymbol{\beta}_t))^2] &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 \cdot 1_{\{y_i=1\}} + p_i(\boldsymbol{\beta}_t)^2 \cdot 1_{\{y_i=0\}}] \\ &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 E[1_{\{y_i=1\}} \mid x_i] + p_i(\boldsymbol{\beta}_t)^2 E[1_{\{y_i=0\}} \mid x_i]] \\ &= E[(1 - p_i(\boldsymbol{\beta}_t))^2 p_i(\boldsymbol{\beta}_t) + p_i(\boldsymbol{\beta}_t)^2 (1 - p_i(\boldsymbol{\beta}_t))] \\ &= E[p_i(\boldsymbol{\beta}_t)(1 - p_i(\boldsymbol{\beta}_t))] \end{aligned}$$

- Achtung: Hier ist \mathbf{X} wieder deterministisch!

- Äquivalente Modellierung:

$$\mathbf{y} = \mathbf{p}(\beta_t) + \boldsymbol{\epsilon}$$

wobei $\boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \in \mathbb{R}^n$ für, bedingt auf \mathbf{X} unabhängige Bernoulli-Fehler, ϵ_j mit

$$E[\boldsymbol{\epsilon} | \mathbf{X}] = \mathbf{0}, \text{Var}[\boldsymbol{\epsilon} | \mathbf{X}] = \mathbf{V} = \begin{pmatrix} v_1 & & \\ & \ddots & \\ & & v_n \end{pmatrix}.$$

$$\mathbf{y} = \mathbf{p}(\beta_t) + \boldsymbol{\epsilon}$$

- Ersetze $\mathbf{y} - \mathbf{p}(\beta_t)$ durch $\boldsymbol{\epsilon}$:

$$i(\beta) = \mathbf{X}^\top \left(\boldsymbol{\epsilon} - \mathbf{V}\mathbf{X}(\beta - \beta_t) - \frac{1}{2} \cdot \begin{pmatrix} (\beta - \beta_t)^\top \ddot{\mathbf{p}}_1(\beta - \beta_t) \\ \vdots \\ (\beta - \beta_t)^\top \ddot{\mathbf{p}}_n(\beta - \beta_t) \end{pmatrix} \right)$$

- Für den MLE $\hat{\beta}$, eingesetzt ist die Score-Funktion $\dot{l}(\cdot)$ gilt:

$$\mathbf{0} = \mathbf{X}^\top \left(\boldsymbol{\epsilon} - \mathbf{VX}(\hat{\beta} - \beta_t) - \frac{1}{2} \cdot \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1(\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n(\hat{\beta} - \beta_t) \end{pmatrix} \right)$$

- Für den MLE $\hat{\beta}$, eingesetzt ist die Score-Funktion $\dot{l}(\cdot)$ gilt:

$$\mathbf{0} = \mathbf{X}^\top \left(\boldsymbol{\epsilon} - \mathbf{VX}(\hat{\beta} - \beta_t) - \frac{1}{2} \cdot \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1(\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n(\hat{\beta} - \beta_t) \end{pmatrix} \right)$$

- Löse nach $(\hat{\beta} - \beta_t)$ auf:

$$\begin{aligned} \mathbf{X}^\top \mathbf{VX}(\hat{\beta} - \beta_t) &= \mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1(\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n(\hat{\beta} - \beta_t) \end{pmatrix} \\ \Leftrightarrow \hat{\beta} - \beta_t &= (\mathbf{X}^\top \mathbf{VX})^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1(\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n(\hat{\beta} - \beta_t) \end{pmatrix} \right) \end{aligned}$$

$$\hat{\beta} - \beta_t = (\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1 (\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n (\hat{\beta} - \beta_t) \end{pmatrix} \right)$$

- Für $\hat{\beta} \rightarrow \beta_t$ ist $\hat{\beta} - \beta_t = \underbrace{(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top}_{=: \mathbf{M}} \boldsymbol{\epsilon}$.

$$\hat{\beta} - \beta_t = (\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_1 (\hat{\beta} - \beta_t) \\ \vdots \\ (\hat{\beta} - \beta_t)^\top \ddot{\mathbf{p}}_n (\hat{\beta} - \beta_t) \end{pmatrix} \right)$$

- Für $\hat{\beta} \rightarrow \beta_t$ ist $\hat{\beta} - \beta_t = \underbrace{(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top}_{=: \mathbf{M}} \boldsymbol{\epsilon}$.
- Wähle $(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}$ als initiale Lösung und substituiere $\hat{\beta} - \beta_t$ auf der rechten Seite:

$$\hat{\beta} - \beta_t \approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} (\mathbf{M}^{-1})^\top \ddot{\mathbf{p}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} (\mathbf{M}^{-1})^\top \ddot{\mathbf{p}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)$$

$$\begin{aligned}\hat{\beta} - \beta_t &\approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} (\mathbf{M}^{-1})^\top \ddot{\mathbf{p}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} (\mathbf{M}^{-1})^\top \ddot{\mathbf{p}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right) \\ &= \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)\end{aligned}$$

- Bemerkung: \mathbf{M}^{-1} ist symmetrisch, da \mathbf{M} symmetrisch ist.

$$\begin{aligned}\mathbf{M}(\mathbf{M}^{-1})^\top &= \mathbf{M}^\top (\mathbf{M}^{-1})^\top = \mathbf{I}_p \\ \Rightarrow \mathbf{M}^{-1} &= (\mathbf{M}^{-1})^\top.\end{aligned}$$

$$\hat{\beta} - \beta_t \approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)$$

- Linker Summand:

$$\begin{aligned} E[(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}] &= E[E[(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \mid \mathbf{X}]] \\ &= E[(\mathbf{X}^\top \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^\top \underbrace{E[\boldsymbol{\epsilon} \mid \mathbf{X}]}_{=0}] \\ &= \mathbf{0} \end{aligned}$$

$$\hat{\beta} - \beta_t \approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\mathbf{p}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)$$

- Rechter Summand: Für \mathbf{X} -messbares \mathbf{A} gilt:

$$\begin{aligned} E[\boldsymbol{\epsilon}^\top \mathbf{A} \boldsymbol{\epsilon} \mid \mathbf{X}] &= E \left[\sum_{i,j=1}^n \mathbf{A}_{ij} \epsilon_i \epsilon_j \mid \mathbf{X} \right] \\ &= \sum_{i,j=1}^n \mathbf{A}_{ij} \cdot \underbrace{E[\epsilon_i \epsilon_j \mid \mathbf{X}]}_{= \text{Cov}(\epsilon_i, \epsilon_j \mid \mathbf{X})} \\ &= \sum_{i=1}^n \mathbf{A}_{ii} \cdot \text{Var}[\epsilon_i \mid \mathbf{X}] \\ &= \text{tr}(\mathbf{A} \text{Var}[\boldsymbol{\epsilon} \mid \mathbf{X}]) \\ &= \text{tr}(\mathbf{A} \mathbf{V}) \end{aligned}$$

$$\hat{\beta} - \beta_t \approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)$$

- Rechter Summand:

$$E \left[\begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \mid \mathbf{X} \right] = \begin{pmatrix} \text{tr}(\mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_1 \mathbf{M}^{-1} \mathbf{X}^\top \mathbf{V}) \\ \vdots \\ \text{tr}(\mathbf{X} \mathbf{M}^{-1} \ddot{\boldsymbol{\rho}}_n \mathbf{M}^{-1} \mathbf{X}^\top \mathbf{V}) \end{pmatrix}$$

$$\ddot{\mathbf{p}}_k = v_k(1 - 2p_k) \cdot x_i x_i^\top$$

- Betrachte zunächst $\text{tr}(\mathbf{X}\mathbf{M}^{-1}\ddot{\mathbf{p}}_k\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V})$ ohne $v_k(1 - 2p_k)$:

$$\begin{aligned} x_k^\top \mathbf{M}^{-1} \mathbf{X}^\top &= x_k^\top (\mathbf{M}^{-1} x_1 \quad \dots \quad \mathbf{M}^{-1} x_n) \\ &= (x_k^\top \mathbf{M}^{-1} x_1 \quad \dots \quad x_k^\top \mathbf{M}^{-1} x_n) \in \mathbb{R}^{1 \times n} \end{aligned}$$

$$\begin{aligned} \text{tr}(\mathbf{X}\mathbf{M}^{-1}x_kx_k^\top\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) &= \text{tr}((x_k^\top\mathbf{M}^{-1}\mathbf{X}^\top)^\top(x_k^\top\mathbf{M}^{-1}\mathbf{X}^\top)\mathbf{V}) \\ &= \sum_{i=1}^n (x_k^\top\mathbf{M}^{-1}x_i)^2 \cdot v_i \\ &= \sum_{i=1}^n (x_k^\top\mathbf{M}^{-1}x_i)(x_i^\top\mathbf{M}^{-1}x_k) \cdot v_i \\ &= x_k^\top\mathbf{M}^{-1} \underbrace{\left(\sum_{i=1}^n v_i \cdot x_i x_i^\top \right)}_{=\mathbf{X}^\top\mathbf{V}\mathbf{X}=\mathbf{M}} \mathbf{M}^{-1}x_k \\ &= x_k^\top\mathbf{M}^{-1}x_k \end{aligned}$$

$$\text{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{x}_k\mathbf{x}_k^\top\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) = \mathbf{x}_k^\top\mathbf{M}^{-1}\mathbf{x}_k$$

- Konstanten lassen sich einfach aus der Spur herausziehen:

$$\begin{aligned} \begin{pmatrix} \text{tr}(\mathbf{X}\mathbf{M}^{-1}\ddot{\mathbf{p}}_1\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) \\ \vdots \\ \text{tr}(\mathbf{X}\mathbf{M}^{-1}\ddot{\mathbf{p}}_n\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) \end{pmatrix} &= \begin{pmatrix} v_1(1 - 2p_1) \text{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{x}_1\mathbf{x}_1^\top\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) \\ \vdots \\ v_n(1 - 2p_n) \text{tr}(\mathbf{X}\mathbf{M}^{-1}\mathbf{x}_n\mathbf{x}_n^\top\mathbf{M}^{-1}\mathbf{X}^\top\mathbf{V}) \end{pmatrix} \\ &= \mathbf{V} \begin{pmatrix} (1 - 2p_1)\mathbf{x}_1^\top\mathbf{M}^{-1}\mathbf{x}_1 \\ \vdots \\ (1 - 2p_n)\mathbf{x}_n^\top\mathbf{M}^{-1}\mathbf{x}_n \end{pmatrix} \end{aligned}$$

$$\hat{\beta} - \beta_t \approx \mathbf{M}^{-1} \left(\mathbf{X}^\top \boldsymbol{\epsilon} - \frac{1}{2} \cdot \mathbf{X}^\top \begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{p}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{p}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \right)$$

- Wir fassen zusammen:

$$E[\hat{\beta} - \beta_t | \mathbf{X}]$$

$$\begin{aligned} &\approx \underbrace{E[\mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} | \mathbf{X}]}_{=0} - \frac{1}{2} \cdot \mathbf{M}^{-1} \mathbf{X}^\top E \left[\begin{pmatrix} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{p}_1 \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \ddot{p}_n \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} | \mathbf{X} \right] \\ &= -\frac{1}{2} \cdot \mathbf{M}^{-1} \mathbf{X}^\top \mathbf{V} \begin{pmatrix} (1 - 2p_1) \mathbf{x}_1^\top \mathbf{M}^{-1} \mathbf{x}_1 \\ \vdots \\ (1 - 2p_n) \mathbf{x}_n^\top \mathbf{M}^{-1} \mathbf{x}_n \end{pmatrix} \end{aligned}$$

- Für den Bias von $\hat{\beta}$ integrieren wir schließlich $E[\hat{\beta} - \beta_t | \mathbf{X}]$ aus.

- In der Praxis: Um den hergeleiteten approximativen Bias korrigieren nicht möglich, da β_t unbekannt ist.

- In der Praxis: Um den hergeleiteten approximativen Bias korrigieren nicht möglich, da β_t unbekannt ist.
- Stattdessen: Bestimme $\hat{\beta}$ und schätze damit \mathbf{V} und \mathbf{p} .

$$\widehat{\text{bias}(\hat{\beta})} = -\frac{1}{2} \cdot \mathbf{M}^{-1} \mathbf{X}^T \hat{\mathbf{V}} \begin{pmatrix} (1 - 2\hat{p}_1) \mathbf{x}_1^T \mathbf{M}^{-1} \mathbf{x}_1 \\ \vdots \\ (1 - 2\hat{p}_n) \mathbf{x}_n^T \mathbf{M}^{-1} \mathbf{x}_n \end{pmatrix}$$

- In der Praxis: Um den hergeleiteten approximativen Bias korrigieren nicht möglich, da β_t unbekannt ist.
- Stattdessen: Bestimme $\hat{\beta}$ und schätze damit \mathbf{V} und \mathbf{p} .

$$\widehat{\text{bias}(\hat{\beta})} = -\frac{1}{2} \cdot \mathbf{M}^{-1} \mathbf{X}^T \hat{\mathbf{V}} \begin{pmatrix} (1 - 2\hat{p}_1) x_1^T \mathbf{M}^{-1} x_1 \\ \vdots \\ (1 - 2\hat{p}_n) x_n^T \mathbf{M}^{-1} x_n \end{pmatrix}$$

- Bias Corrected MLE:

$$\hat{\beta}_{\text{corrected}} = \hat{\beta} - \widehat{\text{bias}(\hat{\beta})}$$

- Wie verhält sich unser (idealer) Bias Correction Term asymptotisch?

The bias is

$$\begin{aligned} \text{bias}(\hat{\beta}) \approx & (-1/2) (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \{ \text{trace}(\mathbf{X} (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} E^k (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\boldsymbol{\varepsilon})) \} \\ & (-1/2) (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V} \{ (1 - 2\pi_k) \mathbf{x}_k^T (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{x}_k \} \end{aligned} \quad (6)$$

where $\{a_k\}$ represents a vector whose k th element is a_k .

By writing a typical element of (6) as a sum and noting that, since we assumed $\lim (\mathbf{X}^T \mathbf{V} \mathbf{X})/n$ finite, the elements of $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$ are $O(1/n)^6$, one can show that the elements in (6) are $O(1/n)$.

- Wie verhält sich unser (idealer) Bias Correction Term asymptotisch?

The bias is

$$\text{bias}(\hat{\beta}) \approx (-1/2) (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \{ \text{trace}(\mathbf{X} (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} E^k (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\boldsymbol{\varepsilon})) \} \\ (-1/2) (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V} \{ (1 - 2\pi_k) \mathbf{x}_k^T (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{x}_k \} \quad (6)$$

where $\{a_k\}$ represents a vector whose k th element is a_k .

By writing a typical element of (6) as a sum and noting that, since we assumed $\lim (\mathbf{X}^T \mathbf{V} \mathbf{X})/n$ finite, the elements of $(\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1}$ are $O(1/n)^6$, one can show that the elements in (6) are $O(1/n)$.

- Wir zeigen:

Proposition

$n \cdot E[\hat{\beta} - \beta_t \mid \mathbf{X}] \xrightarrow{\mathbb{P}} -\frac{1}{2} E[\mathbf{v} \cdot \mathbf{x} \mathbf{x}^T]^{-1} E[(1 - 2p) \mathbf{v} \cdot \mathbf{x} \mathbf{x}^T E[\mathbf{v} \cdot \mathbf{x} \mathbf{x}^T]^{-1} \mathbf{x}] \in \mathbb{R}^p$.
Insbesondere ist $E[\hat{\beta} - \beta_t \mid \mathbf{X}] \in O_p(1/n)$.

Lemma

Die Abbildung $(A \mapsto A^{-1}) : \text{GL}_n(\mathbb{R}) \rightarrow \text{GL}_n(\mathbb{R})$ ist stetig bezüglich des von $(\mathbb{R}^{n \times n}, \|\cdot\|)$ induzierten metrischen Raumes.

Beweis: Invertiere $A \in \text{GL}_n(\mathbb{R})$ mithilfe der Cramerschen Regel

$$A^{-1} = \frac{1}{\det A} ((-1)^{i+j} \det A_{ji})_{ij}$$

wobei A_{ji} die Streichmatrix bezeichnet, welche durch Weglassen der i -ten Zeile und j -ten Spalte entsteht. Dass \det stetig ist, sieht man z.B. anhand der Leibnizformel:

$$\det A = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{i=1}^n a_{i\sigma(i)}.$$



Proposition

$n \cdot E[\hat{\beta} - \beta_t \mid \mathbf{X}] \xrightarrow{\mathbb{P}} -\frac{1}{2} E[v \cdot \mathbf{xx}^\top]^{-1} E[v(1-2p) \cdot \mathbf{xx}^\top E[v \cdot \mathbf{xx}^\top]^{-1} \mathbf{x}] \in \mathbb{R}^p$.
Insbesondere ist $E[\hat{\beta} - \beta_t \mid \mathbf{X}] \in O_p(1/n)$.

Beweis: Da $E[v \cdot \mathbf{xx}^\top]$ endlich ist, folgt mit dem starken Gesetz der Großen Zahlen

$$\frac{1}{n} M = \frac{1}{n} \sum_{i=1}^n v_i \cdot x_i x_i^\top \xrightarrow{f.s.} E[v \cdot \mathbf{xx}^\top].$$

Insbesondere existiert $E[v \cdot \mathbf{xx}^\top]^{-1}$ wegen der positiven Definitheit.
Continuous Mapping Theorem:

$$nM^{-1} = \left(\frac{1}{n} M \right)^{-1} \xrightarrow{f.s.} E[v \cdot \mathbf{xx}^\top]^{-1}.$$

$$\begin{aligned}
n \cdot E[\hat{\beta} - \beta_t | \mathbf{X}] &= -\frac{1}{2} n \cdot M^{-1} \mathbf{X}^T \mathbf{V} \begin{pmatrix} (1 - 2p_1) \mathbf{x}_1^T M^{-1} \mathbf{x}_1 \\ \vdots \\ (1 - 2p_n) \mathbf{x}_n^T M^{-1} \mathbf{x}_n \end{pmatrix} \\
&= -\frac{1}{2} n M^{-1} \sum_{i=1}^n v_i (1 - 2p_i) \mathbf{x}_i \mathbf{x}_i^T M^{-1} \mathbf{x}_i \\
&= -\frac{1}{2} \underbrace{n M^{-1}}_{\xrightarrow{f.s.} E[v \cdot \mathbf{x} \mathbf{x}^T]} \cdot \frac{1}{n} \sum_{i=1}^n v_i (1 - 2p_i) \mathbf{x}_i \mathbf{x}_i^T n M^{-1} \mathbf{x}_i
\end{aligned}$$

Wegen Slutsky's Lemma reicht es also zu zeigen, dass

$$\frac{1}{n} \sum_{i=1}^n v_i (1 - 2p_i) \mathbf{x}_i \mathbf{x}_i^T n M^{-1} \mathbf{x}_i \xrightarrow{\mathbb{P}} E \left[v(1 - 2p) \cdot \mathbf{x} \mathbf{x}^T E[v \cdot \mathbf{x} \mathbf{x}^T]^{-1} \mathbf{x} \right].$$

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n v_i (1 - 2p_i) x_i x_i^\top n M^{-1} x_i \\ &= \frac{1}{n} \sum_{i=1}^n v_i (1 - 2p_i) x_i x_i^\top \left((n M^{-1} - E[v \cdot x x^\top]^{-1}) + E[v \cdot x x^\top]^{-1} \right) x_i \end{aligned}$$

Da $\|x\| \leq C < \infty$ fast sicher, existieren alle Momente $E[\|x\|^n]$ für $n \in \mathbb{N}$.
Wegen der Äquivalenz der Normen, gilt für geeignetes $d > 0$:

$$E \left[\left\| v(1 - 2p) x x^\top E[v \cdot x x^\top]^{-1} x \right\| \right] \leq E \left[d \|x\|^3 \left\| E[v \cdot x x^\top]^{-1} \right\| \right] < \infty.$$

Nach dem schwachen Gesetz der großen Zahlen:

$$\frac{1}{n} \sum_{i=1}^n v_i (1 - 2p_i) x_i x_i^\top E[v \cdot x x^\top]^{-1} x_i \xrightarrow{\mathbb{P}} E \left[v(1 - 2p) \cdot x x^\top E[v \cdot x x^\top]^{-1} x \right].$$

Weiter gilt fast sicher:

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n v_i(1 - 2p_i)x_i x_i^\top (nM^{-1} - E[v \cdot xx^\top]^{-1})x_i \right\| \\ & \leq \frac{1}{n} \sum_{i=1}^n dC^3 \left\| nM^{-1} - E[v \cdot xx^\top]^{-1} \right\| \\ & = dC \left\| nM^{-1} - E[v \cdot xx^\top]^{-1} \right\| \xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

Insgesamt also:

$$\frac{1}{n} \sum_{i=1}^n v_i(1 - 2p_i)x_i x_i^\top nM^{-1}x_i \xrightarrow{\mathbb{P}} E \left[v(1 - 2p) \cdot xx^\top E[v \cdot xx^\top]^{-1}x \right].$$



- Bemerkung: Falls wir wissen, dass $\sup_{n \in \mathbb{N}} \|nM^{-1}\| < \infty$, können wir mit dominierter Konvergenz auch zeigen, dass

$$E[n \cdot E[\hat{\beta} - \beta_t | \mathbf{X}]]$$

$$\xrightarrow{n \rightarrow \infty} -\frac{1}{2} E[v \cdot \mathbf{xx}^\top]^{-1} E \left[v(1 - 2p) \cdot \mathbf{xx}^\top E[v \cdot \mathbf{xx}^\top]^{-1} \mathbf{x} \right]$$

- Bemerkung: Falls wir wissen, dass $\sup_{n \in \mathbb{N}} \|nM^{-1}\| < \infty$, können wir mit dominierter Konvergenz auch zeigen, dass

$$E[n \cdot E[\hat{\beta} - \beta_t | \mathbf{X}]] \\ \xrightarrow{n \rightarrow \infty} -\frac{1}{2} E[v \cdot \mathbf{xx}^\top]^{-1} E \left[v(1 - 2p) \cdot \mathbf{xx}^\top E[v \cdot \mathbf{xx}^\top]^{-1} \mathbf{x} \right]$$

- Hierfür zeigen wir, dass für jede Teilfolge der Erwartungswerte eine konvergente Teilteilfolge existiert.

- Bemerkung: Falls wir wissen, dass $\sup_{n \in \mathbb{N}} \|nM^{-1}\| < \infty$, können wir mit dominierten Konvergenz auch zeigen, dass

$$E[n \cdot E[\hat{\beta} - \beta_t | \mathbf{X}]] \\ \xrightarrow{n \rightarrow \infty} -\frac{1}{2} E[v \cdot \mathbf{xx}^\top]^{-1} E \left[v(1 - 2p) \cdot \mathbf{xx}^\top E[v \cdot \mathbf{xx}^\top]^{-1} \mathbf{x} \right]$$

- Hierfür zeigen wir, dass für jede Teilfolge der Erwartungswerte eine konvergente Teilteilfolge existiert.
- Verwende dazu die äquivalente Formulierung für stochastisch konvergente Zufallsvariablen:

Proposition

$X_n \xrightarrow{\mathbb{P}} X$ genau dann, wenn für jede Teilfolge $(X_{n_m})_{m \in \mathbb{N}}$ eine Teilteilfolge $(X_{n_{m_k}})_{k \in \mathbb{N}}$ existiert, sodass $X_{n_{m_k}} \xrightarrow{f.s.} X$.

- Was passiert mit dem $o\left(\|\hat{\beta} - \beta_t\|^2\right)$ Term?

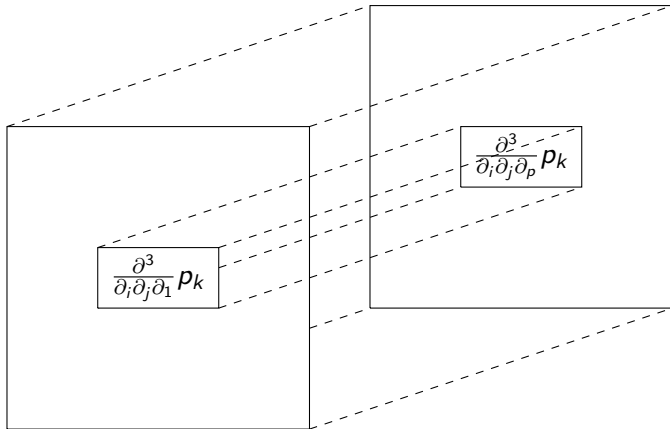
To ensure accounting for all terms of $O(1/n)$, we carried the Taylor series expansion in (3) to the third order. The resulting additional terms were all $O(1/n^2)$, hence (6) represents the bias of $\hat{\beta}$ to order $O(1/n)$.

- Ausrechnen...

- Wir berechnen für jedes $k = 1, \dots, n$ die dritte Ableitung $\ddot{\mathbf{p}}_k$:

$$\ddot{\mathbf{p}}_k = \begin{array}{c} \begin{array}{cc} \frac{\partial^3}{\partial_1 \partial_1 \partial_1} \mathbf{p}_k & \cdots & \frac{\partial^3}{\partial_1 \partial_\rho \partial_1} \mathbf{p}_k \\ \vdots & & \vdots \\ \frac{\partial^3}{\partial_\rho \partial_1 \partial_1} \mathbf{p}_k & \cdots & \frac{\partial^3}{\partial_\rho \partial_\rho \partial_1} \mathbf{p}_k \end{array} \\ \begin{array}{cc} \frac{\partial^3}{\partial_1 \partial_1 \partial_\rho} \mathbf{p}_k & \cdots & \frac{\partial^3}{\partial_1 \partial_\rho \partial_\rho} \mathbf{p}_k \\ \vdots & & \vdots \\ \frac{\partial^3}{\partial_\rho \partial_1 \partial_\rho} \mathbf{p}_k & \cdots & \frac{\partial^3}{\partial_\rho \partial_\rho \partial_\rho} \mathbf{p}_k \end{array} \end{array}$$

$$\ddot{\mathbf{p}}_k^{(i,j)} =$$



$$\begin{aligned}
 &= \frac{d}{d\beta} v_k (1 - 2p_k) x_k^{(i)} x_k^{(j)} \\
 &= \underbrace{(v_k (1 - 2p_k)^2 - 2v_k^2)}_{=: s_k} x_k^\top x_k^{(i)} x_k^{(j)}
 \end{aligned}$$

$$\ddot{\mathbf{p}}_k^{(i,j)} = s_k \mathbf{x}_k^\top \mathbf{x}_k^{(i)} \mathbf{x}_k^{(j)}$$

- Multiplizieren des Tensors $\ddot{\mathbf{p}}_k$ mit den Richtungen $\hat{\beta} - \beta_t$:

$$\frac{1}{6} s_k \left\{ \mathbf{x}_k^{(d)} (\hat{\beta} - \beta_t)^\top \mathbf{x}_k \mathbf{x}_k^\top (\hat{\beta} - \beta_t) \right\}_d^\top (\hat{\beta} - \beta_t)$$

- Einsetzen der Initiallösung $\mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}$ liefert als Term 3. Ordnung:

$$\begin{aligned} & \mathbf{M}^{-1} \sum_{i=1}^n o(\|\hat{\beta} - \beta_t\|^2) x_i \\ &= \frac{1}{6} \cdot \mathbf{M}^{-1} \mathbf{X}^\top \begin{pmatrix} s_1 \left\{ \mathbf{x}_1^{(d)} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \mathbf{x}_1 \mathbf{x}_1^\top \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \right\}_d^\top \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \\ \vdots \\ s_n \left\{ \mathbf{x}_n^{(d)} \boldsymbol{\epsilon}^\top \mathbf{X} \mathbf{M}^{-1} \mathbf{x}_n \mathbf{x}_n^\top \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \right\}_d^\top \mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon} \end{pmatrix} \end{aligned}$$

- Das sieht ganz ähnlich aus, wie der Term 2. Ordnung mit zusätzlichem $O_p(1/n)$ Term $\mathbf{M}^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}$.

- Simulationen für Stichprobengrößen $n \in \{25, 200\}$ und Anzahl der Features $p \in \{2, 5\}$.

- Simulationen für Stichprobengrößen $n \in \{25, 200\}$ und Anzahl der Features $p \in \{2, 5\}$.
- Auf $[0, 1]^p$ uniform verteilte x_1, \dots, x_n .

- Simulationen für Stichprobengrößen $n \in \{25, 200\}$ und Anzahl der Features $p \in \{2, 5\}$.
- Auf $[0, 1]^p$ uniform verteilte x_1, \dots, x_n .
- Wahrer Parameter β_t gewählt als Eigenvektor des kleinsten bzw. größten Eigenwerts von $\mathbf{X}^\top \mathbf{X}$.

- Simulationen für Stichprobengrößen $n \in \{25, 200\}$ und Anzahl der Features $p \in \{2, 5\}$.
- Auf $[0, 1]^p$ uniform verteilte x_1, \dots, x_n .
- Wahrer Parameter β_t gewählt als Eigenvektor des kleinsten bzw. größten Eigenwerts von $\mathbf{X}^\top \mathbf{X}$.
- Bestimme Wahrscheinlichkeiten p_i und ziehe y_i wie folgt:

$$y_i = \begin{cases} 1 & p_i \leq u(0, 1) \\ 0 & \text{sonst} \end{cases}$$

- Simulationen für Stichprobengrößen $n \in \{25, 200\}$ und Anzahl der Features $p \in \{2, 5\}$.
- Auf $[0, 1]^p$ uniform verteilte x_1, \dots, x_n .
- Wahrer Parameter β_t gewählt als Eigenvektor des kleinsten bzw. größten Eigenwerts von $\mathbf{X}^\top \mathbf{X}$.
- Bestimme Wahrscheinlichkeiten p_i und ziehe y_i wie folgt:

$$y_i = \begin{cases} 1 & p_i \leq u(0, 1) \\ 0 & \text{sonst} \end{cases}$$

- Pro Konfiguration 500 Replikationen der obigen Prozedur.

Eigenvalue	$p - 1$	Sample Size	ML	Corrected ML
smallest	2	25	10.38 (12.22)	7.46 (7.94)
		200	4.79 (2.37)	4.68 (2.30)
smallest	5	25	39.75 (74.36)	13.03 (9.99)
		200	5.63 (2.63)	5.38 (2.48)
largest	2	25	16.73 (33.54)	8.91 (7.04)
		200	4.73 (1.05)	4.63 (1.01)
largest	5	25	50.22 (66.25)	12.48 (8.10)
		200	5.94 (1.53)	5.65 (1.41)

Table: Mean-Squared Error für den MLE und den korrigierten MLE. In Klammern jeweils dessen Standardabweichung.

- Approximation des Bias durch Taylorentwicklung der Score-Funktion bis zur zweiten Ordnung.

- Approximation des Bias durch Taylorentwicklung der Score-Funktion bis zur zweiten Ordnung.
- Korrektur des MLE um den approximierten Bias.

- Approximation des Bias durch Taylorentwicklung der Score-Funktion bis zur zweiten Ordnung.
- Korrektur des MLE um den approximierten Bias.
- Approximierter Bias liegt in $O(1/n)$.

- Approximation des Bias durch Taylorentwicklung der Score-Funktion bis zur zweiten Ordnung.
- Korrektur des MLE um den approximierten Bias.
- Approximierter Bias liegt in $O(1/n)$.
- Bei extrem kleinen Datensätzen können substantielle Verbesserung beobachtet werden.

- Approximation des Bias durch Taylorentwicklung der Score-Funktion bis zur zweiten Ordnung.
- Korrektur des MLE um den approximierten Bias.
- Approximierter Bias liegt in $O(1/n)$.
- Bei extrem kleinen Datensätzen können substantielle Verbesserung beobachtet werden.
- Tradeoff: Laufzeit vs. Verbesserung.

- Schaefer, R. (1983). Bias Correction in Maximum Likelihood Logistic Regression. *Statistics in Medicine, Vol.2*
- Rohde, A. (2021). Wahrscheinlichkeitstheorie
- Huber-Klawitter, A. (2019). Lineare Algebra I.