

**week11: Effizienz der Maximum-Likelihood-Schätzer, Teil2:  
Mit Cramer-Rao Abschätzung** (nicht klausurrelevant)

In einer Stochastik II Vorlesung wird gelegentlich die sogenannte Cramer-Rao Abschätzung bewiesen, das ist die folgende Sache: Wir haben Zufallszahlen oder zufällige Größen  $x_1, x_2, \dots, x_n$ , die von einer Wahrscheinlichkeitsverteilung

$$p_{\theta}(x) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) \quad (1)$$

generiert worden sind. Es gelte also

$$\int_{\mathbb{R}^n} p_{\theta}(x) d^n x = \int_{\mathbb{R}^n} p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) dx_1 dx_2 \dots dx_n = 1 \quad (2)$$

Die theta's  $\theta_1, \dots, \theta_m$  sind die Modellparameter. Wenn wir das  $p_{\theta}(x)$  nur als Funktion von  $\theta$  auffassen, weil wir für die  $x_1, \dots, x_n$  die uns gegebenen realisierten Daten einsetzen, dann ist das genau die Likelihood-Funktion,

$$L(\theta) = L(\theta_1, \dots, \theta_m) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) \quad (3)$$

Schreiben wir das  $p_{\theta}(x)$  für unsere Situation, Lineare Regression als statistisches Problem, eben hin: Die Modellparameter sind die Regressionskoeffizienten  $\beta_0, \dots, \beta_p$  und die Standardabweichung  $\sigma$  der normalverteilten Fehler  $\varepsilon_1, \dots, \varepsilon_n$ , wenn wir für unsere gegebenen Daten  $\vec{y}$  und  $\vec{x}_0, \dots, \vec{x}_p$  den Regressionsansatz

$$\vec{y} = \beta_0 \vec{x}_0 + \dots + \beta_p \vec{x}_p + \vec{\varepsilon} \quad (4)$$

oder in Koordinaten

$$y_i = \beta_0 x_{i,0} + \dots + \beta_p x_{i,p} + \varepsilon_i =: \mu_i + \varepsilon_i \quad (5)$$

machen. Die zufälligen Größen, die wir oben mit  $x_1, \dots, x_n$  bezeichnet haben, sind hier jetzt die  $y_1, \dots, y_n$ , wir haben hier also eine W'keitsdichte oder Likelihood-Funktion

$$p_{\beta, \sigma}(y) = p_{\beta_0, \dots, \beta_p, \sigma}(y_1, \dots, y_n) = L(\beta_0, \dots, \beta_p, \sigma) = L(\beta, \sigma) \quad (6)$$

Die hatten wir in dem week7.pdf hergeleitet, das war da die Gleichung (9):

$$L(\beta, \sigma) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right] \right\} \quad (7)$$

Schreiben wir jetzt die Cramer-Rao Abschätzung hin, das war etwa das Theorem 5.3.2 aus dem week10a.pdf aus der Stochastik II Vorlesung aus dem SS2021. Für die Formulierung des Theorems benutzen wir zunächst wieder die Notation von ganz oben, die zufälligen Größen seien also ganz allgemeine  $x_1, \dots, x_n$  mit W'keitsdichte  $p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n)$ .

**Theorem (Cramer-Rao Abschätzung, Stochastik II):** Gegeben sei eine Wahrscheinlichkeitsverteilung  $p_\theta(x) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n)$  mit

$$\int_{\mathbb{R}^n} p_\theta(x) d^n x = \int_{\mathbb{R}^n} p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) dx_1 dx_2 \cdots dx_n = 1 \quad (8)$$

Für  $k = 1, \dots, m$  seien

$$\hat{\theta} = \begin{pmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_m \end{pmatrix} : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad (9)$$

erwartungstreue Schätzer, d.h. es gilt

$$\mathbb{E}[\hat{\theta}_k] := \int_{\mathbb{R}^n} \hat{\theta}_k(x) p_\theta(x) d^n x = \theta_k \quad (10)$$

Weiter sei

$$\text{Cov}(\hat{\theta}) := \left( \text{Cov}[\hat{\theta}_k, \hat{\theta}_\ell] \right)_{k, \ell=1, \dots, m} \in \mathbb{R}^{m \times m} \quad (11)$$

die Covarianz-Matrix von  $\hat{\theta}$ , insbesondere ist also  $\mathbb{V}[\hat{\theta}_k] = \text{Cov}(\hat{\theta})_{k,k}$  die Varianz des  $k$ -ten Schätzers. Wir definieren die sogenannte Fisher-Informationsmatrix  $I(\theta)$  durch

$$I(\theta) := \left( -\mathbb{E} \left[ \frac{\partial^2 \log p_\theta}{\partial \theta_k \partial \theta_\ell} \right] \right)_{k, \ell=1, \dots, m} \in \mathbb{R}^{m \times m} \quad (12)$$

mit

$$\mathbb{E} \left[ \frac{\partial^2 \log p_\theta}{\partial \theta_k \partial \theta_\ell} \right] = \int_{\mathbb{R}^n} \frac{\partial^2 \log p_\theta}{\partial \theta_k \partial \theta_\ell}(x_1, \dots, x_n) p_\theta(x_1, \dots, x_n) d^n x \quad (13)$$

und  $I^{-1}(\theta)$  sei das Inverse von  $I(\theta)$ . Dann gilt:

$$\langle v, \text{Cov}(\hat{\theta}) v \rangle \geq \langle v, I^{-1}(\theta) v \rangle \quad \forall v \in \mathbb{R}^m \quad (14)$$

Insbesondere gilt also, wenn wir für  $v$  den  $k$ -ten Standardbasisvektor nehmen,

$$\mathbb{V}[\hat{\theta}_k] \geq [I^{-1}(\theta)]_{k,k} \quad (15)$$

für jeden erwartungstreuen Schätzer  $\hat{\theta}_k : \mathbb{R}^n \rightarrow \mathbb{R}$ .

Wir wollen jetzt das  $I(\theta)$  und das  $I^{-1}(\theta)$  für unsere Situation, mit der W'keitsdichte gegeben durch (7),

$$p_\theta(x) \rightarrow p_{\beta, \sigma}(y) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right] \right\} \quad (16)$$

berechnen. Damit bekommen wir dann eine untere Schranke für die Varianz von erwartungstreuen Schätzern für die Regressionskoeffizienten. Wenn diese untere Schranke dann identisch ist mit der Varianz der Maximum-Likelihood-Schätzer, und das ist der Fall, heisst das dann also, dass die Maximum-Likelihood-Schätzer minimale Varianz haben, also effizient sind. Im Gegensatz zu den Herleitungen von letzter Woche müssen wir diesmal keine spezielle funktionale Form der betrachteten Schätzer voraussetzen, letztes Mal hatten wir ja eine Menge  $\mathcal{L}$  von linearen (in  $\vec{y}$ ) Schätzern betrachtet. Diese Einschränkung ist jetzt nicht mehr nötig, die einzige Voraussetzung ist die Erwartungstreue, wir betrachten also Schätzer

$$\tilde{\beta}_k : \mathbb{R}^n \rightarrow \mathbb{R} \quad (17)$$

mit

$$\mathbb{E}[\tilde{\beta}_k] = \beta_k \quad (18)$$

Wir müssen die Grössen (13) berechnen. Wir haben mit  $\nu := \sigma^2$

$$p_{\beta,\nu}(y) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\nu}} \exp\left[-\frac{(y_i - \mu_i)^2}{2\nu}\right] \right\} \quad (19)$$

$$\log p_{\beta,\nu}(y) = -\frac{n}{2} \log(2\pi\nu) - \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{2\nu} \quad (20)$$

wobei

$$\mu_i = \mu_i(\beta) = \beta_0 x_{i,0} + \cdots + \beta_p x_{i,p} \quad (21)$$

Also bekommen wir

$$\frac{\partial \log p_{\beta,\nu}}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \log p_{\beta,\nu}}{\partial \mu_i} \frac{\partial \mu_i}{\partial \beta_j} = \sum_{i=1}^n \frac{y_i - \mu_i}{\nu} x_{i,j} \quad (22)$$

und

$$\frac{\partial \log p_{\beta,\nu}}{\partial \nu} = -\frac{n}{2\nu} + \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{2\nu^2} \quad (23)$$

Damit erhalten wir die folgenden zweiten Ableitungen:

$$\frac{\partial^2 \log p_{\beta,\nu}}{\partial \beta_j \partial \beta_k} = \sum_{i=1}^n \frac{-x_{i,k}}{\nu} x_{i,j} = -\frac{\vec{x}_j \vec{x}_k}{\nu} \quad (24)$$

$$\frac{\partial^2 \log p_{\beta,\nu}}{\partial \beta_j \partial \nu} = -\sum_{i=1}^n \frac{y_i - \mu_i}{\nu^2} x_{i,j} \quad (25)$$

$$\frac{\partial^2 \log p_{\beta,\nu}}{\partial \nu^2} = \frac{n}{2\nu^2} - \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{\nu^3} \quad (26)$$

Wegen

$$\mathbb{E}[y_i - \mu_i] = \mathbb{E}[\varepsilon_i] = 0 \quad (27)$$

$$\mathbb{E}[(y_i - \mu_i)^2] = \mathbb{E}[\varepsilon_i^2] = \sigma^2 = \nu \quad (28)$$

bekommen wir dann die folgenden Erwartungswerte:

$$\mathbb{E} \left[ \frac{\partial^2 \log p_{\beta, \nu}}{\partial \beta_j \partial \beta_k} \right] = - \frac{\vec{x}_j \vec{x}_k}{\nu} \quad (29)$$

$$\mathbb{E} \left[ \frac{\partial^2 \log p_{\beta, \nu}}{\partial \beta_j \partial \nu} \right] = 0 \quad (30)$$

$$\mathbb{E} \left[ \frac{\partial^2 \log p_{\beta, \nu}}{\partial \nu^2} \right] = \frac{n}{2\nu^2} - \sum_{i=1}^n \frac{\nu}{\nu^3} = \frac{n}{2\nu^2} - \frac{n}{\nu^2} = - \frac{n}{2\nu^2} \quad (31)$$

Also,

$$I(\theta) = I(\beta, \nu) = \begin{pmatrix} \frac{1}{\nu} (\vec{x}_j \vec{x}_k)_{j,k=0}^p & 0 \\ 0 & \frac{n}{2\nu^2} \end{pmatrix} = \begin{pmatrix} \frac{1}{\nu} X^T X & 0 \\ 0 & \frac{n}{2\nu^2} \end{pmatrix} \quad (32)$$

und damit

$$I^{-1}(\beta, \nu) = \begin{pmatrix} \nu (X^T X)^{-1} & 0 \\ 0 & \frac{2\nu^2}{n} \end{pmatrix} \quad (33)$$

Die Cramer-Rao Abschätzung liefert also in diesem Fall (mit  $j = 0, 1, \dots, p$ )

$$\mathbb{V}[\tilde{\beta}_j] \geq I^{-1}(\beta, \nu)_{j,j} = \nu [(X^T X)^{-1}]_{j,j} = \sigma^2 [(X^T X)^{-1}]_{j,j} \quad (34)$$

$$\mathbb{V}[\tilde{\nu}] \geq I^{-1}(\beta, \nu)_{p+1,p+1} = \frac{2\nu^2}{n} = \frac{2\sigma^4}{n} \quad (35)$$

für beliebige erwartungstreue Schätzer  $\tilde{\beta}_j$  und  $\tilde{\nu} = \tilde{\sigma}^2$ . Der Ausdruck auf der rechten Seite von (34) ist aber gerade die Varianz des Maximum-Likelihood-Schätzers für den  $j$ -ten Regressionskoeffizienten,

$$\mathbb{V}[\hat{\beta}_{\text{ML},j}] = \sigma^2 [(X^T X)^{-1}]_{j,j} \quad (36)$$

das hatten wir ja in dem week9.pdf in dem Theorem 9.1, Teil (a), gezeigt. Also haben die Maximum-Likelihood-Schätzer für die Regressionskoeffizienten minimale Varianz in der Menge aller erwartungstreuen Schätzer, sie sind also effizient. Für die Varianz des  $\hat{s}^2$ , die erwartungstreue Version des Maximum-Likelihood-Schätzers für das  $\sigma^2$ , hatten wir ebenfalls im Theorem 9.1, dann im Teil (b), den Ausdruck

$$\mathbb{V}[\hat{s}^2] = \frac{2\sigma^4}{n - (p + 1)} \quad (37)$$

hergeleitet. Für festes  $n$  ist das ein kleines bisschen grösser als die rechte Seite von (35). Hier könnte man dann zunächst nur sagen, dass das  $\hat{s}^2$  asymptotisch, im Limes  $n \rightarrow \infty$ , effizient ist.