

### week7: Maximum Likelihood Schätzung der Regressionskoeffizienten

Wir wollen von jetzt an das lineare Regressionsproblem nicht mehr als ein deterministisches Minimierungsproblem betrachten, sondern als ein statistisches Problem. Dazu betrachten wir das folgende Setup: Gegeben seien wieder die Datenvektoren

$$\vec{x}_0 = \begin{pmatrix} x_{1,0} \\ x_{2,0} \\ \vdots \\ x_{n,0} \end{pmatrix}, \quad \vec{x}_1 = \begin{pmatrix} x_{1,1} \\ x_{2,1} \\ \vdots \\ x_{n,1} \end{pmatrix}, \quad \dots, \quad \vec{x}_p = \begin{pmatrix} x_{1,p} \\ x_{2,p} \\ \vdots \\ x_{n,p} \end{pmatrix} \quad (1)$$

und

$$\vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad (2)$$

und  $\vec{x}_0$  sei etwa wieder der konstante Vektor mit lauter Einsen,  $\vec{x}_0 = (1, \dots, 1)$  oder  $x_{i,0} = 1$  für  $1 \leq i \leq n$ . Jetzt machen wir die folgende

**Annahme (Regression als statistisches Problem):** Die Ypsilon's  $\vec{y} = (y_1, \dots, y_n)$  sind nicht einfach gegebene Daten, sondern "in Wirklichkeit" kommen sie von einem stochastischen Modell

$$\vec{y} = \beta_0 \vec{x}_0 + \beta_1 \vec{x}_1 + \dots + \beta_p \vec{x}_p + \vec{\varepsilon} \quad (3)$$

her, wobei die  $\vec{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$  normalverteilte, unabhängige Zufallszahlen mit Mittelwert 0 und Standardabweichung  $\sigma$  sind,

$$\varepsilon_i \in N(0, \sigma) \quad \text{unabhängig} \quad (4)$$

Die  $\vec{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$  sind gegebene Zahlen, die wir aber nicht kennen und bestimmen wollen.

**Problem:** Welche Auswahl von

$$\vec{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$$

passt 'am besten' zum statistischen Modell (3,4), wenn die  $\vec{x}_j$  und  $\vec{y}$  konkret durch die Daten (1,2) gegeben sind?

**Lösung:** Wir bestimmen die  $\vec{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$  mit Hilfe der Maximum-Likelihood-Methode, dadurch ist dann auch die genaue mathematische Bedeutung von ‘am besten’ definiert: Die beste Auswahl der  $(\beta_0, \beta_1, \dots, \beta_p)$  soll die Likelihood-Funktion maximieren. Wir müssen also zunächst die Likelihood-Funktion bestimmen:

Jedes einzelne  $y_i$  ist gegeben durch

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \varepsilon_i \\ &= \mu_i + \varepsilon_i \end{aligned} \quad (5)$$

wobei wir abgekürzt haben

$$\mu_i := \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} \quad (6)$$

Insbesondere ist jedes  $y_i$  also selber eine normalverteilte Zufallszahl mit Mittelwert  $\mu_i$  und Standardabweichung  $\sigma$ . Das bedeutet, dass die Wahrscheinlichkeit, dass sich ein solches  $y_i$  in einem Intervall  $[\tilde{y}_i, \tilde{y}_i + d\tilde{y}_i)$  realisiert, gegeben ist durch

$$\text{Prob}[y_i \in [\tilde{y}_i, \tilde{y}_i + d\tilde{y}_i)] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tilde{y}_i - \mu_i)^2}{2\sigma^2}} d\tilde{y}_i \quad (7)$$

Wenn die  $\varepsilon_i$  unabhängig sind, sind auch die  $y_i$  unabhängig und die W'keit, dass sich alle  $y_i$  in gegebenen Intervallen  $[\tilde{y}_i, \tilde{y}_i + d\tilde{y}_i)$  realisieren tun, ist dann also gegeben durch das Produkt

$$\begin{aligned} \text{Prob}[y_1 \in [\tilde{y}_1, \tilde{y}_1 + d\tilde{y}_1), y_2 \in [\tilde{y}_2, \tilde{y}_2 + d\tilde{y}_2), \dots, y_n \in [\tilde{y}_n, \tilde{y}_n + d\tilde{y}_n)] &= \\ &= \text{Prob}[y_1 \in [\tilde{y}_1, \tilde{y}_1 + d\tilde{y}_1)] \text{Prob}[y_2 \in [\tilde{y}_2, \tilde{y}_2 + d\tilde{y}_2)] \dots \text{Prob}[y_n \in [\tilde{y}_n, \tilde{y}_n + d\tilde{y}_n)] \\ &\stackrel{(7)}{=} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tilde{y}_i - \mu_i)^2}{2\sigma^2}} d\tilde{y}_i \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (\tilde{y}_i - \mu_i)^2\right\} d\tilde{y}_1 \dots d\tilde{y}_n \end{aligned} \quad (8)$$

Nun haben wir ja konkret realisierte Werte für die  $y_1, \dots, y_n$ , die wir dann also auf der rechten Seite von (8) für die  $\tilde{y}_i$  einsetzen können. Wenn wir das machen, bekommen wir eine Funktion  $L = L(\vec{\beta}, \sigma)$ , die nur noch von den Modellparametern  $\beta_0, \dots, \beta_p$ , die sind in den  $\mu_i$ 's drin, und  $\sigma$  (und den Intervallbreiten  $d\tilde{y}_i \equiv dy_i$ , die aber wieder aus der eigentlichen Rechnung rausfallen werden) abhängen tut, das ist dann die Likelihood Funktion:

$$L(\vec{\beta}, \sigma) := \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - \mu_i)^2}{2\sigma^2}\right] dy_i \right\} \quad (9)$$

Wir betrachten wieder den Logarithmus (mit den  $\mu_i$  gegeben durch die Gleichung (6) von oben) und bekommen

$$\log L(\vec{\beta}, \sigma) = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2 + \text{const} \quad (10)$$

wobei die Konstante

$$\text{const} := -\frac{n}{2} \log(2\pi) + \sum_{i=1}^n \log(dy_i)$$

nur Terme enthält, die nicht von den Modellparametern  $\vec{\beta}$  und  $\sigma$  abhängen. Das Maximieren von  $L$  ist äquivalent zum Maximieren von  $\log L$ , wir maximieren  $\log L$ . Notwendige Bedingung ist

$$\left( \frac{\partial \log L}{\partial \beta_0}, \frac{\partial \log L}{\partial \beta_1}, \dots, \frac{\partial \log L}{\partial \beta_p}, \frac{\partial \log L}{\partial \sigma} \right) = (0, 0, \dots, 0, 0) \quad (11)$$

Wir haben

$$\begin{aligned} \frac{\partial \log L}{\partial \beta_j} &= + \frac{1}{2\sigma^2} \sum_{i=1}^n 2(y_i - \mu_i) \frac{\partial \mu_i}{\partial \beta_j} \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu_i) x_{i,j} \\ &= \frac{1}{\sigma^2} \left\{ \vec{y} \cdot \vec{x}_j - [\beta_0 \vec{x}_0 + \dots + \beta_p \vec{x}_p] \cdot \vec{x}_j \right\} \stackrel{!}{=} 0 \end{aligned}$$

für alle  $0 \leq j \leq p$ . Damit bekommen wir das Gleichungssystem

$$\begin{aligned} \beta_0 \vec{x}_0 \vec{x}_0 + \beta_1 \vec{x}_0 \vec{x}_1 + \dots + \beta_p \vec{x}_0 \vec{x}_p &= \vec{x}_0 \vec{y} \\ \beta_0 \vec{x}_1 \vec{x}_0 + \beta_1 \vec{x}_1 \vec{x}_1 + \dots + \beta_p \vec{x}_1 \vec{x}_p &= \vec{x}_1 \vec{y} \\ &\vdots \\ \beta_0 \vec{x}_p \vec{x}_0 + \beta_1 \vec{x}_p \vec{x}_1 + \dots + \beta_p \vec{x}_p \vec{x}_p &= \vec{x}_p \vec{y} \end{aligned}$$

oder in Matrix-Notation

$$A \vec{\beta} = \vec{b} \quad (12)$$

mit

$$A = \begin{pmatrix} \vec{x}_0 \vec{x}_0 & \vec{x}_0 \vec{x}_1 & \dots & \vec{x}_0 \vec{x}_p \\ \vec{x}_1 \vec{x}_0 & \vec{x}_1 \vec{x}_1 & \dots & \vec{x}_1 \vec{x}_p \\ \vdots & \vdots & & \vdots \\ \vec{x}_p \vec{x}_0 & \vec{x}_p \vec{x}_1 & \dots & \vec{x}_p \vec{x}_p \end{pmatrix}, \quad \vec{b} = \begin{pmatrix} \vec{x}_0 \vec{y} \\ \vec{x}_1 \vec{y} \\ \vdots \\ \vec{x}_p \vec{y} \end{pmatrix} \quad (13)$$

Das ist exakt dasselbe System, was wir auch schon in `week3.pdf` in den Gleichungen (10) und (11) erhalten hatten. Also bekommen wir auch exakt dieselbe Formel für die Regressionskoeffizienten: Mit der Matrix der Regressoren

$$X := \begin{pmatrix} | & | & & | \\ \vec{x}_0 & \vec{x}_1 & \dots & \vec{x}_p \\ | & | & & | \end{pmatrix} \in \mathbb{R}^{n \times (p+1)} \quad (14)$$

sind die beta's gegeben durch

$$\vec{\beta} = A^{-1} \vec{b} = (X^T X)^{-1} X^T \vec{y} =: \hat{\beta}_{\text{ML}} \quad (15)$$

$\hat{\beta}_{\text{ML}}$  ist also der Maximum-Likelihood-Schätzer für die Regressionskoeffizienten. Das  $\sigma$ , die Standardabweichung der Fehler, können wir auch schätzen: Wir haben

$$\begin{aligned} \frac{\partial \log L}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{2}{2\sigma^3} \sum_{i=1}^n (y_i - \mu_i)^2 \stackrel{!}{=} 0 \\ \Leftrightarrow &\sum_{i=1}^n (y_i - \mu_i)^2 = n \sigma^2 \end{aligned}$$

oder

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mu_i)^2 =: \hat{\sigma}_{\text{ML}}^2 \quad (16)$$

mit

$$\mu_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_p x_{i,p}$$

oder

$$\vec{\mu} = X \vec{\beta} = X (X^T X)^{-1} X^T \vec{y} = P_X \vec{y} \quad (17)$$

Damit können wir auch schreiben

$$\begin{aligned} \hat{\sigma}_{\text{ML}}^2 &= \frac{1}{n} (\vec{y} - \vec{\mu})^2 \\ &= \frac{1}{n} (\vec{y} - P_X \vec{y})^2 \\ &= \frac{1}{n} [(Id - P_X) \vec{y}]^2 \\ &= \frac{1}{n} [P_{X^\perp} \vec{y}]^2 \end{aligned} \quad (18)$$

Fassen wir unsere Ergebnisse in einem Theorem zusammen:

**Theorem 7.1:** Wir betrachten die Formulierung des linearen Regressionsproblems als statistisches Problem gegeben durch die Gleichungen (3) und (4). Dann ist die Likelihood-Funktion  $L = L(\vec{\beta}, \sigma)$  gegeben durch

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

mit

$$\mu_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p}$$

und wird maximiert von den Maximum-Likelihood-Schätzern

$$\hat{\beta}_{\text{ML}} = (X^T X)^{-1} X^T \vec{y}, \quad (19)$$

$$\hat{\sigma}_{\text{ML}}^2 = \frac{1}{n} [P_{X^\perp} \vec{y}]^2. \quad (20)$$

Im folgenden wollen wir die statistischen Eigenschaften der Maximum-Likelihood-Schätzer etwas genauer untersuchen. Dazu schauen wir uns zunächst die Erwartungstreue an. Wir wollen also überprüfen, ob die Gleichungen

$$\mathbf{E}[\hat{\beta}_{\text{ML}}] \stackrel{?}{=} \vec{\beta} \quad (21)$$

$$\mathbf{E}[\hat{\sigma}_{\text{ML}}^2] \stackrel{?}{=} \sigma^2 \quad (22)$$

gültig sind. Dazu müssen wir uns zunächst klarmachen, was das  $E[\cdot]$  hier genau bedeutet, was ist da zu rechnen?

In die Schätzer können wir die  $\vec{x}_0, \dots, \vec{x}_p$  und das  $\vec{y}$  einsetzen und bekommen dann Zahlen heraus. Wir können etwa schreiben

$$\begin{aligned}\hat{\beta}_{\text{ML}} &= \hat{\beta}_{\text{ML}}(X, \vec{y}) : \mathbb{R}^{n \times (p+1)} \times \mathbb{R}^n \rightarrow \mathbb{R}^{p+1} \\ \hat{\sigma}_{\text{ML}}^2 &= \hat{\sigma}_{\text{ML}}^2(X, \vec{y}) : \mathbb{R}^{n \times (p+1)} \times \mathbb{R}^n \rightarrow \mathbb{R}\end{aligned}$$

Für das

$$X = \begin{pmatrix} | & | & \cdots & | \\ \vec{x}_0 & \vec{x}_1 & & \vec{x}_p \\ | & | & & | \end{pmatrix}$$

setzen wir auf jeden Fall unsere konkret gegebenen Daten aus Gleichung (1) ein. Für die  $\vec{y} = (y_1, \dots, y_n)$  können wir entweder die ebenfalls konkret gegebenen Daten aus Gleichung (2) einsetzen, oder wir setzen Daten ein, die von dem stochastischen Modell (3) herkommen, die also durch Zufallszahlen generiert worden sind. Zur Überprüfung der Erwartungstreue machen wir jetzt letzteres, so dass wir etwa schreiben können:

$$\hat{\beta}_{\text{ML}} = \hat{\beta}_{\text{ML}}(X, \vec{y}) = \hat{\beta}_{\text{ML}}(X, X\vec{\beta} + \vec{\varepsilon}) \quad (23)$$

$$\hat{\sigma}_{\text{ML}}^2 = \hat{\sigma}_{\text{ML}}^2(X, \vec{y}) = \hat{\sigma}_{\text{ML}}^2(X, X\vec{\beta} + \vec{\varepsilon}) \quad (24)$$

Dadurch hängen die Schätzer dann also von den  $N(0, \sigma)$ -normalverteilten Zufallszahlen  $\vec{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$  ab. Der Erwartungswert in den Gleichungen (21) und (22) meint dann einfach Integration gegen die entsprechenden  $n$  Gauss'schen Glockenkurven:

$$E[\hat{\beta}_{\text{ML}}] := \int_{\mathbb{R}^n} \hat{\beta}_{\text{ML}}(X, X\vec{\beta} + \vec{\varepsilon}) \prod_{k=1}^n \left\{ e^{-\frac{\varepsilon_k^2}{2\sigma^2}} \frac{d\varepsilon_k}{\sqrt{2\pi\sigma^2}} \right\} \quad (25)$$

$$E[\hat{\sigma}_{\text{ML}}^2] := \int_{\mathbb{R}^n} \hat{\sigma}_{\text{ML}}^2(X, X\vec{\beta} + \vec{\varepsilon}) \prod_{k=1}^n \left\{ e^{-\frac{\varepsilon_k^2}{2\sigma^2}} \frac{d\varepsilon_k}{\sqrt{2\pi\sigma^2}} \right\} \quad (26)$$

Man bekommt dann die folgenden Resultate:

**Theorem 7.2:** a) Der Maximum-Likelihood-Schätzer  $\hat{\beta}_{\text{ML}}$  für die Regressionskoeffizienten ist erwartungstreu, es gilt

$$E[\hat{\beta}_{\text{ML}}] = \vec{\beta} \quad (27)$$

b) Der Maximum-Likelihood-Schätzer  $\hat{\sigma}_{\text{ML}}^2$  für die Varianz der Fehler ist nur asymptotisch, im Limes  $n \rightarrow \infty$  erwartungstreu. Genauer gilt

$$E[\hat{\sigma}_{\text{ML}}^2] = \frac{n-(p+1)}{n} \sigma^2 \quad (28)$$

Damit ist der modifizierte Schätzer

$$\hat{s}^2 := \frac{n}{n-(p+1)} \hat{\sigma}_{\text{ML}}^2 = \frac{1}{n-(p+1)} [P_{X+\vec{y}}]^2 \quad (29)$$

erwartungstreu, es gilt  $E[\hat{s}^2] = \sigma^2$ .

**Beweis:** ..machen wir nächste Woche.