

week11a: Kapitel 7: Effizienz von Schätzern, Teil 1

Die Maximum-Likelihood-Schätzer für den Mittelwert und die Varianz von normalverteilten Zufallszahlen aus dem Beispiel 1 waren gegeben durch

$$\hat{\mu}(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

$$\hat{\sigma}^2(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2)$$

mit $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i = \hat{\mu}$. Das $\hat{\mu}$ war erwartungstreu,

$$\mathbb{E}[\hat{\mu}] = \mu \quad (3)$$

aber das $\hat{\sigma}^2$ war nur im Limes $n \rightarrow \infty$ erwartungstreu. Deshalb hatten wir den erwartungstreuen Schätzer \hat{s}^2 mit Vorfaktor $1/(n-1)$ definiert, das war also

$$\hat{s}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (4)$$

und wir haben dann die Identität

$$\mathbb{E}[\hat{s}^2] = \sigma^2 \quad (5)$$

Wir wollen jetzt zeigen, dass diese Schätzer effizient sind. Das heisst, dass sie innerhalb der Menge aller erwartungstreuen Schätzer eine minimale Varianz besitzen. Das ist jetzt schon eine etwas nichttrivialere Sache, weil das einzige, was man hier fordert, Erwartungstreue ist, ansonsten können die Schätzer, mit denen man das $\hat{\mu}$ und das \hat{s}^2 vergleichen möchte, völlig beliebige Funktionen sein. Also definieren wir diese Menge von Schätzern, mit denen wir das $\hat{\mu}$ und das \hat{s}^2 vergleichen wollen (etwa \mathcal{M} für Mittelwert und \mathcal{V} für Varianz):

$$\mathcal{M} := \left\{ \tilde{\mu} : \mathbb{R}^n \rightarrow \mathbb{R} \mid \mathbb{E}[\tilde{\mu}(x_1, \dots, x_n)] = \mu \right\} \quad (6)$$

$$\mathcal{V} := \left\{ \tilde{\sigma}^2 : \mathbb{R}^n \rightarrow \mathbb{R} \mid \mathbb{E}[\tilde{\sigma}^2(x_1, \dots, x_n)] = \sigma^2 \right\} \quad (7)$$

wobei die Erwartungswerte gemäss

$$\mathbb{E}[F(x_1, \dots, x_n)] = \int_{\mathbb{R}^n} F(x_1, \dots, x_n) \prod_{i=1}^n \left\{ e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \frac{dx_i}{\sqrt{2\pi\sigma^2}} \right\}$$

zu berechnen sind. Dann gilt das folgende

Theorem 7.1.1: a) Der Maximum-Likelihood-Schätzer $\hat{\mu}$ aus Gleichung (1) ist effizient. Das heisst, für jeden Schätzer $\tilde{\mu}$ aus \mathcal{M} gilt:

$$V[\tilde{\mu}] \geq \frac{\sigma^2}{n} = V[\hat{\mu}] \quad \forall \tilde{\mu} \in \mathcal{M} \quad (8)$$

b) Die erwartungstreue Version (4) des Maximum-Likelihood-Schätzers für die Varianz ist asymptotisch effizient. Das heisst genauer, für jeden Schätzer $\tilde{\sigma}^2$ aus \mathcal{V} gilt:

$$V[\tilde{\sigma}^2] \geq \frac{2\sigma^4}{n} = \frac{n-1}{n} V[\hat{s}^2] \stackrel{n \rightarrow \infty}{\approx} V[\hat{s}^2] \quad \forall \tilde{\sigma}^2 \in \mathcal{V} \quad (9)$$

Beweis: ..machen wir gleich mit der Cramer-Rao Abschätzung.

Der Beweis ergibt sich als Folgerung aus der sogenannten Cramer-Rao Abschätzung. Das ist eine untere Schranke für die Varianz von erwartungstreuen Schätzern, und das bemerkenswerte an dieser Schranke ist, dass man sie für ein sehr allgemeines Setting hinschreiben kann. Wir müssen nicht unbedingt eine unabhängige Folge von Zufallszahlen haben wie in den Beispielen 1 und 2, sondern die Cramer-Rao Abschätzung funktioniert auch noch für das Beispiel 3. Das genaue Setting ist das folgende:

Wir haben Zufallszahlen oder zufällige Grössen x_1, x_2, \dots, x_n , die von einer Wahrscheinlichkeitsverteilung

$$p_\theta(x) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) \quad (10)$$

generiert worden sind. Es gelte also

$$\int_{\mathbb{R}^n} p_\theta(x) d^n x = \int_{\mathbb{R}^n} p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) dx_1 dx_2 \dots dx_n = 1 \quad (11)$$

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

$$L(\theta) = L(\theta_1, \dots, \theta_m) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) \quad (12)$$

Schreiben wir das $p_\theta(x)$ für unsere 3 Beispiele noch einmal hin:

Beispiel 1: Für Beispiel 1 haben wir

$$p_\theta(x) = p_{\mu, \sigma^2}(x_1, \dots, x_n) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right\} \quad (13)$$

Beispiel 2: Für Beispiel 2 haben wir

$$p_\theta(x) = p_\lambda(x_1, \dots, x_n) = \prod_{i=1}^n \left\{ \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right\} \quad (14)$$

mit $x_i \in \{0, 1, 2, \dots\}$, die x_i sind hier also diskret. In einem solchen Fall wollen wir in der Gleichung (11) das Integral dann immer als eine Summe interpretieren, also hier wäre das dann

$$\sum_{x_1, \dots, x_n=0}^{\infty} p_\lambda(x_1, \dots, x_n) = 1$$

mit $\theta = \theta_1 = \lambda$ als einzigen Modellparameter.

Beispiel 3: Und für das Beispiel 3, den zeitdiskreten Ornstein-Uhlenbeck Prozess, hatten wir (mit $x_k := x_{t_k}$)

$$p_\theta(x) = p_{\alpha, \beta, \eta}(x_1, \dots, x_n) = \prod_{k=1}^n \left\{ \frac{1}{\sqrt{2\pi\eta^2}} e^{-\frac{(x_k - [\alpha x_{k-1} + \beta])^2}{2\eta^2}} \right\} \quad (15)$$

mit den Modellparametern α, β, η oder den ursprünglichen Modellparametern

$$\begin{aligned} \kappa &= \frac{1-\alpha}{\Delta t} \\ \bar{x} &= \frac{\beta}{1-\alpha} \\ \sigma^2 &= \frac{\eta^2}{\Delta t} \end{aligned}$$

Auch in diesem Fall gilt die Gleichung (11).

In allen 3 Beispielen können wir Ableitungen nach den Modellparametern mit den Integralen oder den Summen über die x_i vertauschen, es gilt also

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta_k} \underbrace{\int_{\mathbb{R}^n} p_\theta(x) d^n x}_{=1} = \frac{\partial}{\partial \theta_k} \int_{\mathbb{R}^n} p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n) dx_1 dx_2 \cdots dx_n \\ &= \int_{\mathbb{R}^n} \frac{\partial p_{\theta_1, \dots, \theta_m}}{\partial \theta_k}(x_1, \dots, x_n) dx_1 dx_2 \cdots dx_n \end{aligned} \quad (16)$$

Dass man die Ableitung mit dem Integral vertauschen kann, ist jetzt nicht immer ganz selbstverständlich, da wir etwa bei gleichverteilten Zufallszahlen ja eine Rechteck-Funktion als W'keitsdichte haben, und die ist im klassischen Sinne nicht differenzierbar (nur im Distributions-Sinn), und solche Fälle wollen wir hier mal ausschliessen. Schreiben wir jetzt die Cramer-Rao Abschätzung hin:

Theorem 7.1.2 (Cramer-Rao Lower Bound): Gegeben sei eine Wahrscheinlichkeitsverteilung $p_\theta(x) = p_{\theta_1, \dots, \theta_m}(x_1, \dots, x_n)$ mit

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

Für $k = 1, \dots, m$ seien

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

erwartungstreue Schätzer, d.h. es gilt

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

Es bezeichne

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

die Covarianz-Matrix von $\hat{\theta}$, insbesondere ist dann also

$$\mathbf{V}[\hat{\theta}_k] = \text{Cov}(\hat{\theta})_{k,k}$$

die Varianz von $\hat{\theta}_k$. Wir definieren die sogenannte Fisher-Informations-Matrix $I(\theta)$ durch

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

mit

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

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

$$\langle v, \text{Cov}(\hat{\theta}) v \rangle \geq \langle v, I^{-1}(\theta) v \rangle \quad \forall v \in \mathbb{R}^m \quad (17)$$

Insbesondere gilt also:

$$\mathbf{V}[\hat{\theta}_k] \geq [I^{-1}(\theta)]_{k,k} \quad (18)$$

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

Bevor wir die Cramer-Rao Abschätzung in der nächsten Veranstaltung beweisen, wollen wir die Abschätzung (18) konkret für das Beispiel 1 hinschreiben und damit das Theorem 7.1.1 vom Anfang der Vorlesung beweisen:

Beweis Theorem 7.1.1: Wir haben mit $\nu := \sigma^2$

$$p_\theta(x) = p_{\mu,\nu}(x_1, \dots, x_n) = \prod_{i=1}^n \left\{ \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{(x_i - \mu)^2}{2\nu}} \right\}$$

$$\log p_\theta(x) = \log p_{\mu,\nu}(x_1, \dots, x_n) = -\frac{n}{2} \log(2\pi\nu) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\nu}$$

und damit

$$\begin{aligned}\frac{\partial \log p_{\mu, \nu}}{\partial \mu} &= \sum_{i=1}^n \frac{x_i - \mu}{\nu} \\ \frac{\partial \log p_{\mu, \nu}}{\partial \nu} &= -\frac{n}{2\nu} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\nu^2}\end{aligned}$$

und

$$\begin{aligned}\frac{\partial^2 \log p_{\mu, \nu}}{\partial \mu^2} &= \sum_{i=1}^n \frac{-1}{\nu} = -\frac{n}{\nu} \\ \frac{\partial^2 \log p_{\mu, \nu}}{\partial \nu \partial \mu} &= -\sum_{i=1}^n \frac{x_i - \mu}{\nu^2} \\ \frac{\partial^2 \log p_{\mu, \nu}}{\partial \nu^2} &= \frac{n}{2\nu^2} - \sum_{i=1}^n \frac{(x_i - \mu)^2}{\nu^3}\end{aligned}$$

Wegen

$$\begin{aligned}\mathbb{E}[x_i - \mu] &= 0 \\ \mathbb{E}[(x_i - \mu)^2] &= \sigma^2 = \nu\end{aligned}$$

bekommen wir dann die folgenden Erwartungswerte:

$$\begin{aligned}\mathbb{E}\left[\frac{\partial^2 \log p_{\mu, \nu}}{\partial \mu^2}\right] &= -\frac{n}{\nu} \\ \mathbb{E}\left[\frac{\partial^2 \log p_{\mu, \nu}}{\partial \nu \partial \mu}\right] &= 0 \\ \mathbb{E}\left[\frac{\partial^2 \log p_{\mu, \nu}}{\partial \nu^2}\right] &= \frac{n}{2\nu^2} - \sum_{i=1}^n \frac{\nu}{\nu^3} = \frac{n}{2\nu^2} - \frac{n}{\nu^2} = -\frac{n}{2\nu^2}\end{aligned}$$

Also,

$$I(\theta) = I(\mu, \nu) = \begin{pmatrix} \frac{n}{\nu} & 0 \\ 0 & \frac{n}{2\nu^2} \end{pmatrix}$$

und damit

$$I^{-1}(\mu, \nu) = \begin{pmatrix} \frac{\nu}{n} & 0 \\ 0 & \frac{2\nu^2}{n} \end{pmatrix}$$

Die Cramer-Rao Abschätzung liefert also in diesem Fall

$$\begin{aligned}\mathbb{V}[\tilde{\mu}] &\geq I^{-1}(\mu, \nu)_{1,1} = \frac{\nu}{n} = \frac{\sigma^2}{n} \\ \mathbb{V}[\tilde{\nu}] &\geq I^{-1}(\mu, \nu)_{2,2} = \frac{2\nu^2}{n} = \frac{2\sigma^4}{n}\end{aligned}$$

für beliebige erwartungstreue Schätzer $\tilde{\mu}$ und $\tilde{\nu} = \tilde{\sigma}^2$ und damit haben wir das Theorem 7.1.1 bewiesen. ■