

Kapitel 6

Ideen und Begriffe aus der Statistik

6.2

6.1 Grundlegende Begriffe, Schätzen von Parametern

Beispiel 6.1. Bei einer biologischen Expedition wurden $n = 53$ Krebse einer gewissen Art gefangen, davon $k = 23$ Weibchen.

Was sagt uns dies über den Weibchenanteil in der Population?

Gibt die Beobachtung (Weibchenanteil in der Stichprobe $\frac{23}{53} \approx 0,434$) Anlass, an einem ausgeglichenen Geschlechterverhältnis in dieser Population zu zweifeln?

Vorstellung: Eine sehr große Population von Krebsen mit (uns unbekanntem) Weibchenanteil $\vartheta \in (0, 1)$, der naheliegendste Schätzwert für ϑ (angesichts der Beobachtungen) ist

$$\widehat{\vartheta} = \frac{23}{53} \quad (\approx 0,434).$$

Bei einer biologischen Expedition wurden $n = 53$ Krebse einer gewissen Art gefangen, davon $k = 23$ Weibchen.

Was sagt uns dies über den Weibchenanteil ϑ in der Population?

Modell: $X = (X_1, X_2, \dots, X_n)$,

$$X_i = \begin{cases} 1, & i\text{-ter gefangener Krebs ist Weibchen,} \\ 0, & i\text{-ter gefangener Krebs ist Männchen} \end{cases}$$

X_i sind u.a., $\sim \text{Ber}_{\vartheta}$ (Wir tun hier so, als ob wir mit Zurücklegen rein zufällig aus der Gesamtpopulation gezogen hätten – diese Approximation ist für große Populationen gerechtfertigt.)

Wir interpretieren $\frac{k}{n} = \frac{23}{53}$ als Realisierung der Zufallsvariable

$$\widehat{\vartheta} := \frac{1}{n}(X_1 + \dots + X_n)$$

Modell:

$$X_i = \begin{cases} 1, & i\text{-ter gefangener Krebs ist Weibchen,} \\ 0, & i\text{-ter gefangener Krebs ist Männchen} \end{cases}$$

X_i sind u.a., $\sim \text{Ber}_\vartheta$, wobei ϑ der (uns unbekannt) tatsächliche Anteil der Weibchen in der Population ist.

i. (Punktschätzung) Was auch immer ϑ ist, es gilt

$$\hat{\vartheta} = \frac{1}{n} (X_1 + \dots + X_n)$$

$$\mathbb{E}_\vartheta[\hat{\vartheta}] = \vartheta,$$

$$\text{Var}_\vartheta[\hat{\vartheta}] = \frac{1}{n^2} n \vartheta (1 - \vartheta) = \frac{1}{n} \vartheta (1 - \vartheta) = \frac{\sigma^2}{n}$$

mit $\sigma = \sigma(\vartheta) = \sqrt{\vartheta(1 - \vartheta)}$. (\mathbb{E}_ϑ , etc. bezieht sich auf Erwartungswerte bezüglich dem Wahrscheinlichkeitsmaß, unter dem $X_i \sim \text{Ber}_\vartheta$ und u.a. sind.)

$$\hat{\sigma} := \sqrt{\hat{\vartheta}(1 - \hat{\vartheta})}$$

ist ein naheliegender Schätzer für σ .

Ein Schätzer für die Standardabweichung von $\hat{\vartheta}$ ist $\frac{\hat{\sigma}}{\sqrt{n}}$.

Im Statistik-Jargon heißt dies auch der „Standardfehler“ (englisch: standard error [of the mean], SEM), dies ist eine naheliegende Maßzahl für die „Genauigkeit der Schätzung“.

(Im Beispiel:

$$\frac{23}{53} \cdot \frac{30}{53} \approx 0,246, \quad \hat{\sigma} \approx 0,496, \quad \frac{\hat{\sigma}}{\sqrt{53}} \approx 0,0681$$

man gibt also an: geschätzter Weibchenanteil $0,43 \pm 0,068$.)

Modell:

$$X_i = \begin{cases} 1, & i\text{-ter gefangener Krebs ist Weibchen,} \\ 0, & i\text{-ter gefangener Krebs ist Männchen} \end{cases}$$

X_i sind u.a., $\sim \text{Ber}_\vartheta$, wobei ϑ der (uns unbekannt) tatsächliche Anteil der Weibchen in der Population ist.

2. („Wie genau ist die Schätzung?“: **Konfidenzintervall**) Wenn der wahre Weibchenanteil ϑ ist, so ist $X_1 + \dots + X_n \sim \text{Bin}_{n,\vartheta}$, also

$$\frac{\hat{\vartheta} - \vartheta}{\hat{\sigma}/\sqrt{n}} \approx \frac{\hat{\vartheta} - \vartheta}{\sigma(\vartheta)/\sqrt{n}} \stackrel{d}{\approx} \mathcal{N}_{0,1}$$

gemäß dem Satz von de Moivre-Laplace (Satz 5.1 und Korollar 5.2), also

$$P_\vartheta\left(-1,96 \leq \frac{\hat{\vartheta} - \vartheta}{\hat{\sigma}/\sqrt{n}} \leq 1,96\right) \approx P(-1,96 \leq Z \leq 1,96) \approx 0,95$$

mit $Z \sim \mathcal{N}_{0,1}$, d.h. das (zufällige) Intervall

$$I := \left[\hat{\vartheta} - 1,96 \frac{\hat{\sigma}}{\sqrt{n}}, \hat{\vartheta} + 1,96 \frac{\hat{\sigma}}{\sqrt{n}}\right] \text{ erfüllt } P_\vartheta(\vartheta \in I) \approx 0,95$$

(für jede Wahl von $\vartheta \in [0, 1]$).

I heißt ein Konfidenzintervall (für ϑ) zum (approximativen) Niveau 0,95.

(Im Beispiel: $I \approx [0,30, 0,57]$)

$$\hat{\vartheta} - 1,96 \cdot \frac{\hat{\sigma}}{\sqrt{n}} \leq \vartheta$$

$$\vdots$$

$$+$$

$$\vdots$$

Modell:

$$X_i = \begin{cases} 1, & i\text{-ter gefangener Krebs ist Weibchen,} \\ 0, & i\text{-ter gefangener Krebs ist Männchen} \end{cases}$$

X_i sind u.a., $\sim \text{Ber}_{\vartheta}$, wobei ϑ der (uns unbekannt) tatsächliche Anteil der Weibchen in der Population ist.

3. (Testen von Hypothesen) Passen die Beobachtungen zur Hypothese, dass der wahre Weibchenanteil in der Population $\vartheta_0 = \frac{1}{2}$ ist?

Wir beobachten

$$|\widehat{\vartheta} - \frac{1}{2}| = \left| \frac{23}{53} - \frac{1}{2} \right| \approx 0,07,$$

es ist

$$P_{\vartheta_0}(|\widehat{\vartheta} - \vartheta_0| \geq 0,07) \approx P(|Z| \geq \sqrt{4 \cdot n} \cdot 0,07) \approx 2(1 - \Phi(0,96)) \approx 0,336.$$

Demnach: Wenn die Hypothese $\vartheta = \vartheta_0 = \frac{1}{2}$ zutrifft, würden wir in ca. 1/3 der Fälle eine mindestens so große Abweichung $|\widehat{\vartheta} - \frac{1}{2}|$ wie die tatsächlich anhand der Daten beobachtete finden. Insoweit gibt die Beobachtung keinen Anlass, diese Hypothese anzuzweifeln.

Definition 6.2. Ein *statistisches Modell* ist ein Tripel $(\mathcal{M} =) (\mathcal{X}, \mathcal{F}, (P_\vartheta)_{\vartheta \in \Theta})$, wo $\mathcal{X} \neq \emptyset$ Menge („Beobachtungs- oder Stichprobenraum“), $\mathcal{F} \subset 2^{\mathcal{X}}$ eine σ -Algebra, Θ eine Menge (mit $|\Theta| > 1$) und für jedes $\vartheta \in \Theta$ ist P_ϑ ein W’maß auf $(\mathcal{X}, \mathcal{F})$.

Das Modell \mathcal{M} heißt *parametrisch*, wenn $\Theta \subset \mathbb{R}^d$ für ein $d \in \mathbb{N}$, speziell *einparametrisch*, wenn $d = 1$.

\mathcal{M} heißt *diskret*, wenn \mathcal{X} abzählbar ist, \mathcal{M} heißt *stetig*, wenn $\mathcal{X} \subset \mathbb{R}^n$ und jedes P_ϑ eine Dichte $\rho_\vartheta : \mathcal{X} \rightarrow [0, \infty]$ besitzt.

Ein diskretes oder stetiges Modell heißt ein *Standardmodell*.

$(\mathcal{X}^n, \mathcal{F}^{\otimes n}, (P_\vartheta^{\otimes n})_{\vartheta \in \Theta})$ heißt das n -fache Produktmodell von \mathcal{M} (für Produktmaße vgl. Def. 2.21).

Definition 6.3. $(\mathcal{X}, \mathcal{F}, (P_\vartheta)_{\vartheta \in \Theta})$ statistisches Modell, (S, \mathcal{A}) messbarer Raum.

1. Eine Zufallsvariable X (definiert auf $(\mathcal{X}, \mathcal{F})$ mit Werten in S , d.h. $X : \mathcal{X} \rightarrow S$ ist \mathcal{F} - \mathcal{A} -messbar) heißt eine *Statistik* (manchmal auch: „Stichprobe“).
2. Sei $\tau : \Theta \rightarrow \mathbb{R}$ eine reelle Kenngröße (oder „Parametermerkmal“), eine Statistik $T : \mathcal{X} \rightarrow \mathbb{R}$ heißt ein *Schätzer* (genauer: „Punktschätzer“) für τ .
3. Ein Schätzer T für τ heißt *erwartungstreu* (oder „unverzerrt“), wenn gilt

$$\forall \vartheta \in \Theta : \mathbb{E}_\vartheta[T] = \tau(\vartheta).$$

$b_\vartheta(T) := \mathbb{E}_\vartheta[T] - \tau(\vartheta)$ heißt die *Verzerrung* (englisch: bias) von T .

Die typische Konstruktion / Situation eines Schätzers ist $T = t(X)$ für eine Funktion $t : S \rightarrow \mathbb{R}$.

Man schreibt / benennt einen Schätzer für τ oft $\hat{\tau}$.