

Vorlesung 10

Schätzen mit Verlass:

Konfidenzintervalle

1. Schätzen von Anteilen

(Buch S. 121-122)

Große Population (♀ und ♂)
mit unbekanntem Weibchenanteil p

In einer Stichprobe vom Umfang $n = 53$
gab es 23 Weibchen.

Wie zuverlässig ist $\frac{23}{53}$ als Schätzung für p ?

Goldene Idee der Statistik:

In einem idealisiert gedachten Szenario
interpretiert man den Schätzwert
als Realisierung einer Zufallsvariablen
und rechnet mit der Variabilität dieser Zufallsvariablen.

In unserem Eingangsbeispiel wird
die Stichprobenziehung (idealisiert!)
als p -Münzwurf gedeutet.

Als *Schätzer* für p betrachten wir die *Zufallsvariable*

$$H := \frac{K}{n},$$

mit $K :=$ Anzahl der “Erfolge”.

H ist die relative Häufigkeit der Erfolge (die “Trefferquote”).

$$\sigma_H = ?$$

K ist $\text{Bin}(n, p)$ -verteilt. Also:

$$\sigma_H = \frac{1}{\sqrt{n}} \sqrt{p(1-p)}$$

Der Zentrale Grenzwertsatz liefert uns:

H ist approximativ normalverteilt

mit Erwartungswert $\mu_H = p$ und Standardabweichung σ_H .

Also insbesondere:

$$\mathbf{P}_p(|p - H| \leq 2\sigma_H) \approx 0.95$$

$$\mathbf{P}_p(p \in [H - 2\sigma_H, H + 2\sigma_H]) \approx 0.95$$

Das zufällige Intervall $[H - 2\sigma_H, H + 2\sigma_H]$ überdeckt den Parameter p mit Wahrscheinlichkeit ≈ 0.95 .

In der Praxis ist auch σ_H
(aus der *einen* vorliegenden Stichprobe) zu schätzen.

Als Schätzer für $\sigma_H = \frac{1}{\sqrt{n}}\sqrt{p(1-p)}$ bietet sich an:

$$\widehat{\sigma}_H := \frac{1}{\sqrt{n}}\sqrt{H(1-H)}$$

$$\mathbf{P}_p(p \in [H - 2\sigma_H, H + 2\sigma_H]) \approx 0.95$$

überträgt sich auf

$$\mathbf{P}_p(p \in [H - 2\widehat{\sigma}_H, H + 2\widehat{\sigma}_H]) \approx 0.95.$$

Das zufällige Intervall

$$I := [H - 2\widehat{\sigma}_H, H + 2\widehat{\sigma}_H]$$

ist ein

Konfidenzintervall für p

mit approximativer Überdeckungswahrscheinlichkeit 0.95

oder kurz ein

approximatives 95%-Konfidenzintervall für p .

In unserem Eingangsbeispiel ($n = 53$, $k = 23$)
hatten die beobachteten Realisierungen von H und $\widehat{\sigma}_H$
die Werte

$$h = 23/53 = 0.43,$$
$$\frac{1}{\sqrt{n}}\sqrt{h(1-h)} = \sqrt{\frac{0.43 \cdot 0.57}{53}} = 0.07.$$

Als Realisierung von $I = [H - 2\widehat{\sigma}_H, H + 2\widehat{\sigma}_H]$ ergab sich
 $[0.43 - 2 \cdot 0.07, 0.43 + 2 \cdot 0.07] = [0.29, 0.57]$.

Man beachte:

Nicht der Parameter p ist zufällig, sondern das Intervall I .

Im Jargon der Statistik wird oft sowohl der Schätzer (die Zufallsvariable) H als auch der Schätzwert (die Zahl) h mit dem Symbol \hat{p} bezeichnet.

Auf der nächsten Folie,
die das Obige zusammenfasst,
steht \hat{p} für die Zufallsvariable H
(vgl. Buch S. 122):

Das zufällige Intervall

$$I := \left[\hat{p} - \frac{2}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})}, \hat{p} + \frac{2}{\sqrt{n}} \sqrt{\hat{p}(1 - \hat{p})} \right]$$

ist ein

Konfidenzintervall für p

mit approximativer Überdeckungswahrscheinlichkeit 0.95

oder kurz ein

approximatives 95%-Konfidenzintervall für p .

Dabei sollte $np(1 - p)$ “nicht zu klein” sein.

Eine Faustregel für die Anwendbarkeit ist: $nh \geq 9$ und $n(1 - h) \geq 9$.

Ein frischer Blick auf Konfidenzintervalle

(samt Wiederholung des Obigen):

n sei fest, und K sei die Anzahl der Treffer bei n Versuchen und Erfolgswahrscheinlichkeit p ,

d.h. es gilt

$$\mathbf{P}_p(K = k) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

Wie oben setzen wir $H := \frac{K}{n}$, $\hat{p} := \frac{k}{n}$.

Im Geist von Übungsaufgabe 27.S definieren wir folgende
Teilmenge von $\{0, 1, \dots, n\} \times [0, 1]$:

$$M_N := \left\{ (\hat{p}, p) : |\hat{p} - p| \leq 2 \cdot \sqrt{\frac{1}{n} \cdot \hat{p} (1 - \hat{p})} \right\}$$

Für nicht zu kleines npq gilt

wegen der approximativen Normalität von K

(und weil $\sqrt{\frac{1}{n}H(1-H)}$ ein passender Schätzer für σ_H ist):

$$\mathbf{P}_p((H, p) \in M_N) \approx 0.95.$$

Für die Intervalle $I_N(\hat{p}) := \{p : (\hat{p}, p) \in M_N\}$ haben wir für jedes $p \in [0, 1]$ die Ereignisgleichheit

$$\{p \in I_N(H)\} = \{(H, p) \in M_N\}.$$

Also gilt für nicht zu kleines npq :

$$\mathbf{P}_p(p \in I_N(H)) \approx 0.95.$$

So lange die Normalapproximation passt, gilt

$$M_N \approx \left\{ \left(\frac{k}{n}, p \right) : k \text{ liegt zwischen dem } 0.025\text{-Quantil und dem } 0.975\text{-Quantil der Bin}(n, p)\text{-Verteilung, } 0 \leq p \leq 1. \right\}$$

Für $p = 0$ bzw. $p = 1$ wird die Approximation unbrauchbar:

$$\begin{aligned} \left(\frac{k}{n}, 0 \right) &\text{ gehört zu } M_N \text{ nur für } k = 0, \\ \text{und } \left(\frac{k}{n}, 1 \right) &\text{ gehört zu } M_N \text{ nur für } k = n. \end{aligned}$$

Auch für kleine npq wird die Approximation schlecht.

Vorschlag von Clopper und Pearson (1934):

$$M_C := \left\{ \binom{k}{n}, p : k \text{ liegt zwischen dem 0.025-Quantil und dem 0.975-Quantil der Bin}(n, p)\text{-Verteilung, } 0 \leq p \leq 1. \right\}$$

Damit bekommen wir für *alle* $p \in [0, 1]$:

$$\mathbf{P}_p((H, p) \in M_C) \geq 0.95.$$

Nach demselben Muster wie oben setzen wir

$$I_C(\hat{p}) := \{p : (\hat{p}, p) \in M_C\}$$

und bekommen jetzt für *alle* $p \in [0, 1]$:

$$\mathbf{P}_p(p \in I_C(H)) \geq 0.95.$$

Beispiel: Von 1500 Senioren bekam die (per Losentscheid ermittelte) eine Hälfte (die *Behandlungsgruppe*) einen Impfstoff verabreicht und die andere Hälfte (die *Kontrollgruppe*) ein Placebo.

Fünf der Senioren haben sich infiziert, alle gehörten zur Kontrollgruppe.

(Quelle:

<https://www.nejm.org/doi/pdf/10.1056/NEJMoa2034577?articleTo>

Table 3, Age group ≥ 75 yr)

Als Realisierung des 95% Clopper-Pearson Konfidenzintervalls für die Erfolgswahrscheinlichkeit des Impfstoffes ergibt sich $[0.478, 1]$ (denn $0.478^5 = 0.025$).

Dieses Intervall schließt den Parameterwert $p = 1/2$ ein: auch wenn alles rein zufällig zugegangen wäre, hätte ein so extremes Ergebnis eine Wahrscheinlichkeit von mehr als 5%.

2. Schätzung des Erwartungswertes einer Verteilung auf \mathbb{R} (Lageschätzung)

$$m := \frac{1}{n}(x_1 + \cdots + x_n)$$

wird gedacht als eine Realisierung der Zufallsvariablen

$$\bar{X} := \frac{1}{n}(X_1 + \cdots + X_n)$$

mit X_1, \dots, X_n unabhängig, identisch verteilt

mit Erwartungswert μ und Standardabweichung σ .

Anders als bei der Anteilschätzung ist hier σ i.a. keine Funktion von μ .

\bar{X} ist ein Schätzer für μ .

\bar{X} hat Erwartungswert μ und Standardabweichung σ/\sqrt{n} .

Ein Schätzer für σ^2 ist

$$S^2 := \frac{1}{n-1} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right)$$

Im Jargon der Statistik schreibt man oft s^2 statt S^2 und verwendet s^2 dann auch zur Bezeichnung von Schätzwerten.

Anders als im Buch werden wir hier auf den Folien an der oben definierten Bezeichnung S^2 festhalten und s^2 für die Bezeichnung einer Realisierung von S^2 reservieren.

2a. Großer Stichprobenumfang n

Der Zentrale Grenzwertsatz liefert uns:

$\bar{X} - \mu$ ist approximativ $N(0, \frac{\sigma^2}{n})$ -verteilt.

Bei bekanntem σ ist also für große n

$$\left[\bar{X} - 2\frac{\sigma}{\sqrt{n}}, \bar{X} + 2\frac{\sigma}{\sqrt{n}} \right]$$

ein approximatives 95%-Konfidenzintervall für μ .

In der Praxis hat man auch σ aus den Daten zu schätzen.

Für große n ist auch

$$J := \left[\bar{X} - 2 \frac{S}{\sqrt{n}}, \bar{X} + 2 \frac{S}{\sqrt{n}} \right]$$

ein approximatives 95%-Konfidenzintervall für μ

Dieses Intervall enthält keine unbekannt Parameter.

Der beobachtete Wert (die Realisierung) von S ist

$$s := \sqrt{\frac{1}{n-1} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)}$$

s ist ein Schätzwert für σ .

Die Realisierung von J ist somit $\left[\bar{x} - 2 \frac{s}{\sqrt{n}}, \bar{x} + 2 \frac{s}{\sqrt{n}} \right]$.

2b. Kleiner Stichprobenumfang n

Für kleine n (etwa: $n \leq 10$) und (exakt bzw. annähernd)
normalverteilte X_i

macht man sich zunutze, dass die Verteilung von

$T := \frac{\sqrt{n}(\bar{X} - \mu)}{S}$ (exakt bzw. annähernd) so verteilt ist wie

$$T_{n-1} := \frac{N_0}{\sqrt{\frac{1}{n-1} (N_1^2 + \dots + N_{n-1}^2)}}$$

mit unabhängigen und $N(0, 1)$ verteilten N_0, \dots, N_{n-1} .

Satz (W. Gosset (alias “Student”, 1908), R. Fisher (1924))

Sind X_1, \dots, X_n unabhängig und $N(\mu, \sigma^2)$ -verteilt,

dann ist $T := \frac{\sqrt{n}(\bar{X} - \mu)}{S}$ so verteilt wie

$$T_{n-1} := \frac{N_0}{\sqrt{\frac{1}{n-1} (N_1^2 + \dots + N_{n-1}^2)}}$$

mit unabhängigen und $N(0, 1)$ verteilten N_0, \dots, N_{n-1} .

Das folgt aus der Rotationssymmetrie
der Standard-Normalverteilung im \mathbb{R}^n

mit einem ähnlichen Argument wie dem in F7a1.12, vgl. Buch S. 138.

Satz (W. Gosset (alias “Student”, 1908), R. Fisher (1924))

Sind X_1, \dots, X_n unabhängig und $N(\mu, \sigma^2)$ -verteilt,

dann ist $T := \frac{\sqrt{n}(\bar{X} - \mu)}{S}$ so verteilt wie

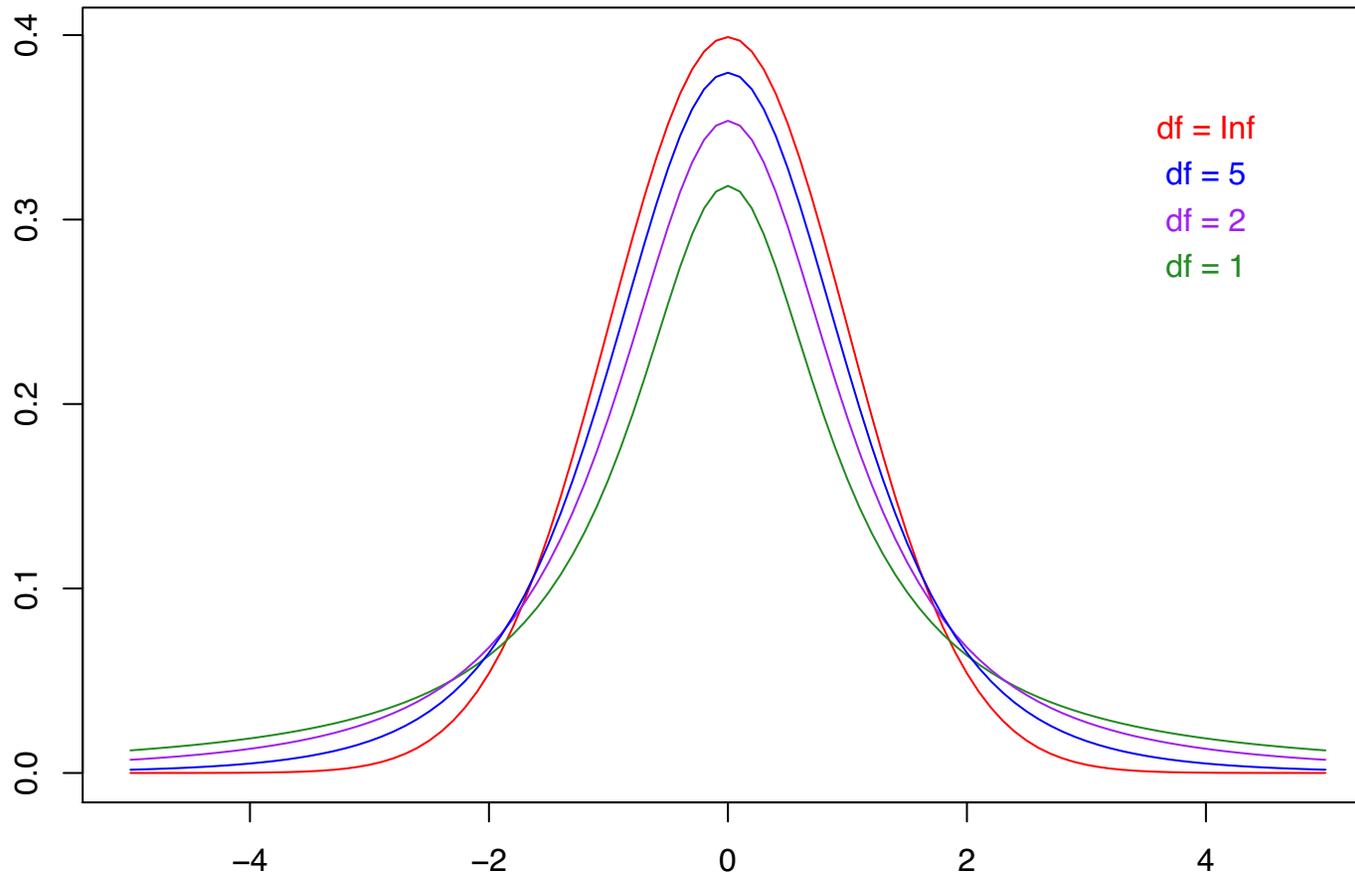
$$T_{n-1} := \frac{N_0}{\sqrt{\frac{1}{n-1} (N_1^2 + \dots + N_{n-1}^2)}}$$

mit unabhängigen und $N(0, 1)$ verteilten N_0, \dots, N_{n-1} .

T heißt (die aus X_1, \dots, X_n gewonnene) **t -Statistik**.

Die **Verteilung von T_{n-1}** heißt **t -Verteilung** (oder **Student-Verteilung**) mit $n - 1$ Freiheitsgraden.

Student's t: Dichtefunktionen



“df” steht hier für “degrees of freedom” , d.h. “Freiheitsgrade”

$$T_{n-1} := \frac{N_0}{\sqrt{\frac{1}{n-1} (N_1^2 + \dots + N_{n-1}^2)}}$$

Für $n \rightarrow \infty$ ist T_{n-1} asymptotisch $N(0, 1)$ -verteilt
(Gesetz der großen Zahlen für den Nenner von T_{n-1}).

Je kleiner n , um so mehr schwankt der Nenner, und
um so *breitschultriger* ist die Verteilung von T_{n-1}

Z.B. für $n = 6$: $\mathbf{P}(|T_5| \leq 2.57) = 0.95$.

Satz (W. Gosset (alias "Student", 1908), R. Fisher (1924))

Sind X_1, \dots, X_n unabhängig und $N(\mu, \sigma^2)$ -verteilt,

dann ist $T := \frac{\sqrt{n}(\bar{X} - \mu)}{S}$ so verteilt wie

$$T_{n-1} := \frac{N_0}{\sqrt{\frac{1}{n-1} (N_1^2 + \dots + N_{n-1}^2)}}$$

mit unabhängigen und $N(0, 1)$ verteilten N_0, \dots, N_{n-1} .

Folgerung: Sind X_1, \dots, X_n unabhängig und $N(\mu, \sigma^2)$ -verteilt, dann ist für jedes $c > 0$:

$$\mathbf{P}(|T_{n-1}| \leq c) = \mathbf{P}\left(\left|\frac{\sqrt{n}(\bar{X} - \mu)}{S}\right| \leq c\right) = \mathbf{P}\left(\mu \in \left[\bar{X} - \frac{cS}{\sqrt{n}}, \bar{X} + \frac{cS}{\sqrt{n}}\right]\right)$$

Für ein 95%-Konfidenzintervall bestimme c so, dass sich **für die linke Seite** 0.95 ergibt.

Z.B. für $n = 6$: $\mathbf{P}(|T_5| \leq 2.57) = 0.95$.

Der passende R-Befehl ist `qt(0.975, 5)`, mit der Ausgabe 2.57

Denn: $\mathbf{P}(T_5 \leq 2.57) = 0.975$.

Man sagt: Das 0.975-Quantil der $t(5)$ -Verteilung ist 2.57.

3. Ein Konfidenzintervall für den Median

(Buch S. 128)

Seien X_1, X_2, \dots, X_n unabhängig, mit Verteilung ρ .

Es gibt Situationen, in denen die Schätzung der “Lage” von ρ über den Stichprobenmittelwert problematisch ist – etwa wenn ρ so viel Masse “weit draußen” hat, dass $\sigma^2 = \text{Var}[X_1]$ sehr groß ist.

In diesem Fall ist es günstig, einen “robusteren” Lageschätzer zu verwenden:

Eine Zahl ν heißt *Median* der Verteilung ρ auf \mathbb{R} ,
wenn $\rho((-\infty, \nu]) \geq 1/2$ **und** $\rho([\nu, \infty)) \geq 1/2$ gilt.

Wenn es nur eine Zahl ν gibt mit $\rho((-\infty, \nu]) = 1/2$
(also z. B. wenn $\rho([\ell, r]) = 1$ gilt
und ρ eine strikt positive Dichtefunktion besitzt),
dann ist ν **der** Median von ρ .

Eine Zahl ν heißt *Median* der Verteilung ρ auf \mathbb{R} ,
wenn $\rho((-\infty, \nu]) \geq 1/2$ **und** $\rho([\nu, \infty)) \geq 1/2$ gilt.

Wenn es mehrere Zahlen ν gibt mit $\rho((-\infty, \nu]) = 1/2$,
dann ist jede dieser Zahlen **ein** Median von ρ .

Bsp: Für die uniforme Verteilung auf $[0, 1] \cup [2, 3]$
ist jedes $\nu \in [1, 2]$ ein Median.

Wie schätzt man den Median?

Die *Ordnungsstatistiken* $X_{(1)} \leq \dots \leq X_{(n)}$
sind die aufsteigend geordneten X_1, \dots, X_n .

Sei j eine natürliche Zahl mit $0 \leq j < n/2$.

Ein Kandidat für ein **Konfidenzintervall für den Median** ist

$$[X_{(1+j)}, X_{(n-j)}].$$

Z.B. für $j = 0$:

$$\mathbf{P}_\rho(\nu \notin [X_{(1)}, X_{(n)}]) = \mathbf{P}_\rho(X_{(1)} > \nu) + \mathbf{P}_\rho(X_{(n)} < \nu) .$$

$$\mathbf{P}_\rho(X_{(1)} > \nu) \leq 2^{-n}$$

$$\mathbf{P}_\rho(X_{(n)} < \nu) \leq 2^{-n} .$$

Also:

$$\mathbf{P}_\rho(\nu \in [X_{(1)}, X_{(n)}]) \geq 1 - \frac{1}{2^{n-1}} .$$

$$\mathbf{P}_\rho(\nu \in [X_{(1)}, X_{(n)}]) \geq 1 - \frac{1}{2^{n-1}}.$$

Bsp: $n = 6$

$$\mathbf{P}_\rho(\nu \in [X_{(1)}, X_{(6)}]) \geq 1 - \frac{1}{32} = 0.97 :$$

Das Konfidenzintervall $[X_{(1)}, X_{(6)}]$ für den Median hält die Überdeckungswahrscheinlichkeit 0.97 ein.