

Wahrscheinlichkeitsrechnung und  
Statistik für Biologen  
**7. Konfidenzintervalle**

Martin Hutzenthaler & Dirk Metzler

1. Juni 2009

# Inhaltsverzeichnis

<b>1</b>	<b>Konfidenzintervalle für Erwartungswerte</b>	<b>1</b>
1.1	Beispiel: Carapaxlänge des Springkrebses . . . . .	1
1.2	Theorie . . . . .	2
<b>2</b>	<b>Konfidenzintervalle für Wahrscheinlichkeiten</b>	<b>3</b>
2.1	Beispiel: Porzellankrebs . . . . .	3
2.2	Theorie . . . . .	4
2.3	Beispiel: Porzellankrebs . . . . .	4
2.4	Beispiel: Stockente . . . . .	5
2.5	Anmerkungen . . . . .	6

## 1 Konfidenzintervalle für Erwartungswerte

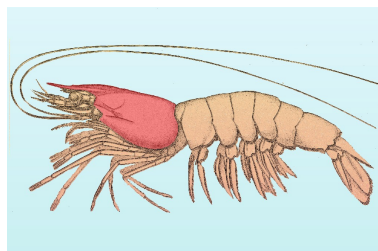
### 1.1 Beispiel: Carapaxlänge des Springkrebses

Beispiel: Springkrebs



*Galathea squamifera*  
image (c) by Matthias Buschmann

Carapaxlänge:



(c): public domain

Wie groß ist die mittlere Carapaxlänge des weiblichen Springkrebses?

Alle weiblichen Springkrebse (also die Grundgesamtheit) zu vermessen, ist zu aufwändig.

Idee: Aus einer Stichprobe läßt sich die mittlere Carapaxlänge schätzen.

Wie genau ist diese Schätzung?

Ziel: Ein Intervall, in dem der Mittelwert der Carapaxlängen aller weiblichen Springkrebse mit hoher Wahrscheinlichkeit liegt.

Dieses Intervall nennen wir **Konfidenzintervall** oder **Vertrauensbereich** für den wahren Wert.

Galathea: Carapaxlänge in einer Stichprobe

Weibchen:  $\bar{x} = 3.23$  mm  $sd(x) = 0.9$  mm  $n = 29$   $sem(x) = \frac{sd(x)}{\sqrt{n}} = \frac{0.9}{\sqrt{29}} = 0.17$  ( $= sd(\bar{x})$ )

Wir kennen bereits folgende Faustformeln:

- **2/3**-Faustformel: Der wahre Mittelwert liegt im Intervall

$$[\bar{x} - sem(x), \bar{x} + sem(x)]$$

mit Wahrscheinlichkeit nahe bei 2/3

- **95%**-Faustformel: Der wahre Mittelwert liegt im Intervall

$$[\bar{x} - 2 * sem(x), \bar{x} + 2 * sem(x)]$$

mit Wahrscheinlichkeit nahe bei 95%

Nun exakt: Sei  $t_{5\%} \leftarrow -qt(0.025, length(x)-1)$ . Dann liegt der wahre Mittelwert mit Wahrscheinlichkeit 95% im Intervall

$$[\bar{x} - t_{5\%} * sem(x), \bar{x} + t_{5\%} * sem(x)]$$

Zur Begründung siehe nächster Abschnitt.

Setzt man die Zahlenwerte  $\bar{x} = 3.23$ ,  $t_{5\%} = 2.05$  und  $sem(x) = 0.17$  in

$$[\bar{x} - t_{5\%} * sem(x), \bar{x} + t_{5\%} * sem(x)]$$

ein, so erhält man das Konfidenzintervall

$$[2.88, 3.58]$$

für den wahren Mittelwert zum Irrtumsniveau 5%.

## 1.2 Theorie

### Konfidenzintervall für den wahren Mittelwert

Ziel: **Bestimme das Konfidenzintervall für den wahren Mittelwert** zum Irrtumsniveau  $\alpha$

Das Konfidenzintervall für den wahren Mittelwert zum Irrtumsniveau  $\alpha$  ist ein aus den Daten  $x = (x_1, \dots, x_n)$  geschätztes (zufälliges) Intervall

$$[\underline{I}(x), \bar{I}(x)]$$

mit folgender Eigenschaft: Ist der wahre Mittelwert gleich  $\mu$  und ist  $(x_1, \dots, x_n)$  eine Stichprobe aus der Grundgesamtheit (mit Mittelwert  $\mu$ ), so gilt

$$\Pr(\mu \in [\underline{I}(x), \bar{I}(x)]) \geq 1 - \alpha$$

Selbstverständlich wollen wir das Konfidenzintervall möglichst klein wählen.

### Konfidenzintervall für den wahren Mittelwert

Lösung: Wir wissen bereits (->Normalapproximation), dass die t-Statistik

$$t := \frac{\bar{x} - \mu}{\text{sem}(x)}$$

annähernd Student-verteilt ist mit  $\text{length}(x)-1$  Freiheitsgraden (wenn  $\text{length}(x)$  groß genug ist).

Sei  $t_\alpha \leftarrow -\text{qt}(\alpha/2, \text{length}(x)-1)$  das  $\alpha/2$ -Quantil der Student-Verteilung mit  $\text{length}(x)-1$  Freiheitsgraden. Dann ist

$$[\bar{x} - t_\alpha * \text{sem}(x), \bar{x} + t_\alpha * \text{sem}(x)]$$

das Konfidenzintervall zum Irrtumsniveau  $\alpha$ .

Begründung:

$$\begin{aligned} & \Pr(\mu \in [\bar{x} - t_\alpha * \text{sem}(x), \bar{x} + t_\alpha * \text{sem}(x)]) \\ &= \Pr(\mu - \bar{x} \in [-t_\alpha * \text{sem}(x), t_\alpha * \text{sem}(x)]) \\ &= \Pr\left(\frac{\mu - \bar{x}}{\text{sem}(x)} \in [-t_\alpha, t_\alpha]\right) \\ &= \Pr\left(\left|\frac{\mu - \bar{x}}{\text{sem}(x)}\right| \leq t_\alpha\right) \\ &= \Pr(|t| \leq t_\alpha) \\ &= 1 - \alpha \end{aligned}$$

Beachte:  $t_\alpha$  wird gerade so gewählt, dass die letzte Gleichung richtig ist.

### Konfidenzintervall allgemein

Sei  $\theta$  ein Parameter der zu Grunde liegenden Verteilung.

Das Konfidenzintervall für den Parameter  $\theta$  zum Irrtumsniveau  $\alpha$  ist ein aus den Daten  $x = (x_1, \dots, x_n)$  geschätztes (zufälliges) Intervall

$$[\underline{I}(x), \bar{I}(x)]$$

mit folgender Eigenschaft: Ist der wahre Parameter gleich  $\theta$  und ist  $(x_1, \dots, x_n)$  eine Stichprobe aus der Grundgesamtheit (mit Parameter  $\theta$ ), so gilt

$$\Pr(\theta \in [\underline{I}(x), \bar{I}(x)]) \geq 1 - \alpha$$

Wie das Konfidenzintervall im Allgemeinen aussieht, ist unbekannt.

## 2 Konfidenzintervalle für Wahrscheinlichkeiten

### 2.1 Beispiel: Porzellankrebs



(c): public domain

Familie: *Porcellanidae*

In einem Fang vom 21.02.1992 in der Helgoländer Tiefe Rinne waren 23 Weibchen und 30 Männchen (*Pisidiae longicornis*), d.h. der Männchenanteil in der Stichprobe war  $30/53 = 0,57$ .

Was sagt uns dies über den Männchenanteil in der Population?

Was ist ein 95%-Konfidenzintervall für den Männchenanteil in der Population? ( $0,57 \pm ??$ )

## 2.2 Theorie

Wir beobachten  $X$  Männchen in einer Stichprobe der Größe  $n$  und möchten den (unbekannten) Männchenanteil  $p$  in der Gesamtpopulation schätzen.

Der offensichtliche Schätzer ist die relative Häufigkeit  $\hat{p} := \frac{X}{n}$  in der Stichprobe.

Frage: Wie verlässlich ist die Schätzung?

Gewünscht: Ein in Abhängigkeit von den Beobachtungen konstruiertes (und möglichst kurzes) Intervall  $[\hat{p}_u, \hat{p}_o]$  mit der Eigenschaft

$$\Pr_p \left( [\hat{p}_u, \hat{p}_o] \text{ überdeckt } p \right) \geq 1 - \alpha$$

für jede Wahl von  $p$ .

Lösungsweg:

Für gegebenes  $p$  ist  $X$  Binomial( $n, p$ )-verteilt,  $E[X] = np$ ,  $\text{Var}[X] = np(1 - p)$ .

Wir wissen: Der Schätzer  $\hat{p}$  ist (in etwa) normalverteilt mit Erwartungswert  $p$  und Standardabweichung  $\sqrt{\frac{p(1-p)}{n}}$ . Dies können wir allerdings nicht verwenden, da  $p$  unbekannt ist.

Stattdessen nutzen wir wieder die Student-Verteilung, wobei hier als Schätzer für die Standardabweichung von  $\hat{p}$

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

verwendet wird. (Das Quadrat hiervon ist ein erwartungstreuer Schätzer der Varianz. In der Literatur wird auch der ML-Schätzer  $\sqrt{\hat{p}(1 - \hat{p})/n}$  verwendet.)

Lösung:

Sei  $\hat{p}$  die relative Häufigkeit in der Stichprobe der Länge  $n$ . Das Konfidenzintervall zum Irrtumsniveau  $\alpha$  ist

$$\left[ \hat{p} - t_{\alpha, n-1} * \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1}}, \hat{p} + t_{\alpha, n-1} * \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1}} \right]$$

wobei  $t_{\alpha, n-1} \leftarrow -\text{qt}(\alpha/2, n-1)$  das  $\alpha/2$ -Quantil der Student-Verteilung mit  $n-1$  Freiheitsgraden ist.

## 2.3 Beispiel: Porzellankrebs

### Männchenanteil beim Porzellankrebs

Setzt man die Zahlenwerte  $n = 53$ ,  $\hat{p} = 0.566$ ,  $t_{5\%, 52} = 2.007$  und  $\sqrt{\hat{p}(1 - \hat{p})/(n - 1)} = 0.0687$  in

$$\left[ \hat{p} - t_{5\%, 52} * \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1}}, \hat{p} + t_{5\%, 52} * \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1}} \right]$$

ein, so erhält man das Konfidenzintervall

$$[0.428, 0.704] = 0.566 \pm 0.138$$

für den wahren Männchenanteil zum Irrtumsniveau 5%.

## 2.4 Beispiel: Stockente



image (c) Andreas Trepte

*Anas platyrhynchos*

Stockente (engl. mallard)

Füchse jagen Stockenten. Durch ihre auffällige Färbung sind dabei Männchen leichter zu erspähen. Hat dies einen Einfluss auf das Geschlechterverhältnis bei amerikanischen Stockenten?

Daten: Stichprobe der Länge  $n = 2200$ . Relative Häufigkeit der Männchen war 0.564.

Daten aus:

## Literatur

[Smi68] Johnson, Sargeant (1977) Impact of red fox predation on the sex ratio of prairie mallards *United States fish & wild life service*

Setzt man die Zahlenwerte  $n = 2200$ ,  $\hat{p} = 0.564$ ,  $t_{5\%,2199} = 1.961$  und  $\sqrt{\hat{p}(1-\hat{p})/(n-1)} = 0.011$  in

$$\left[ \hat{p} - t_{5\%} * \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1}}, \hat{p} + t_{5\%} * \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1}} \right]$$

ein, so erhält man das Konfidenzintervall

$$[0.543, 0.585] = 0.564 \pm 0.021$$

für den wahren Männchenanteil zum Irrtumsniveau 5%.

## 2.5 Anmerkungen

- Für die Gültigkeit der Approximation muss  $n$  genügend groß und  $p$  nicht zu nahe an 0 oder 1 sein. (Eine häufig zitierte „Faustregel“ ist „ $np \geq 9, n(1 - p) \geq 9$ “.)
- Die Philosophie der Konfidenzintervalle entstammt der *frequentistischen* Interpretation der Statistik: Für jede Wahl des Parameters  $p$  würden wir bei häufiger Wiederholung des Experiments finden, dass in (ca.)  $(1 - \alpha) \cdot 100\%$  der Fälle das (zufällige) Konfidenzintervall den „wahren“ (festen) Parameter  $p$  überdeckt.