

# 7. Ideen aus der Statistik

## 7.1 Deskriptive Statistik

Matthias Birkner

[http://www.mathematik.uni-mainz.de/~birkner/GrundlStoch\\_1314/](http://www.mathematik.uni-mainz.de/~birkner/GrundlStoch_1314/)

13.1.2014

# Inhalt

- 1 **Ansatz der Statistik**
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

Viele Menschen stehen „Statistik“ kritisch gegenüber:

*It is easy to lie with statistics.*

*It is easy to lie with statistics.  
It is hard to tell the truth without it.*

Andrejs Dunkels (1939–1998)

# Worum geht es in der Statistik?

Die Natur ist voller Variabilität.

# Worum geht es in der Statistik?

Die Natur ist voller Variabilität.

Wie geht man mit variablen Daten um?

# Worum geht es in der Statistik?

Die Natur ist voller Variabilität.

Wie geht man mit variablen Daten um?

Idee der Statistik:

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

# Worum geht es in der Statistik?

Die Natur ist voller Variabilität.

Wie geht man mit variablen Daten um?

Idee der Statistik:

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

Die Daten werden als Realisierungen von Zufallsvariablen  
aufgefasst, die in einem stochastischen Modell spezifiziert  
werden.



# Worum geht es in der Statistik?

Die Natur ist voller Variabilität.

Wie geht man mit variablen Daten um?

Idee der Statistik:

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

Die Daten werden als Realisierungen von Zufallsvariablen  
aufgefasst, die in einem stochastischen Modell spezifiziert  
werden.

Man versucht dann, anhand der Daten Rückschlüsse auf  
Parameter des Modells zu ziehen, und so systematische Effekte  
von Zufälligem zu trennen.

# Deskriptive (d.h. beschreibende) Statistik

Wie geht man mit variablen Daten um?

# Deskriptive (d.h. beschreibende) Statistik

Wie geht man mit variablen Daten um?

„0. Antwort“: Man verschafft sich einen ersten Eindruck mittels graphischer Darstellungen und statistischer Kenngrößen

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

# Der Springkrebs *Galathea intermedia*



(Daten aus einer Diplomarbeit aus 2001 am Forschungsinstitut Senckenberg, Frankfurt am Main, Crustaceensektion, Leitung: Prof. Michael Türkay)

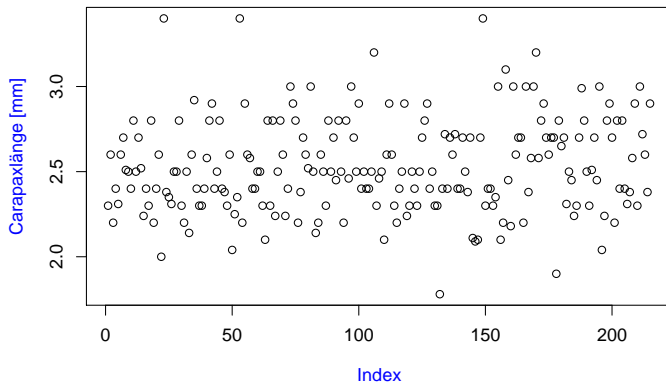
# Helgoländer Tiefe Rinne, Fang vom 6.9.1988

Carapaxlänge (mm):

Nichteiertragende Weibchen ( $n = 215$ )

2,9	3,0	2,9	2,5	2,7	2,9	2,9	3,0
3,0	2,9	3,4	2,8	2,9	2,8	2,8	2,4
2,8	2,5	2,7	3,0	2,9	3,2	3,1	3,0
2,7	2,5	3,0	2,8	2,8	2,8	2,7	3,0
2,6	3,0	2,9	2,8	2,9	2,9	2,3	2,7
2,6	2,7	2,5	.	.	.	.	.

## Nichteiertragende Weibchen am 6. Sept. '88, n=215



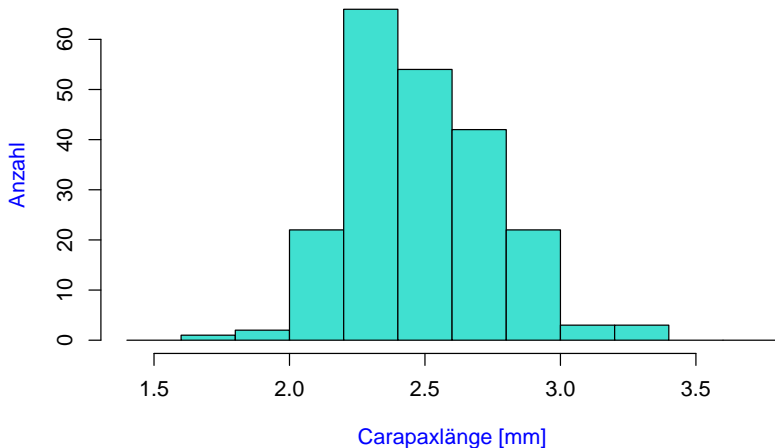
# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

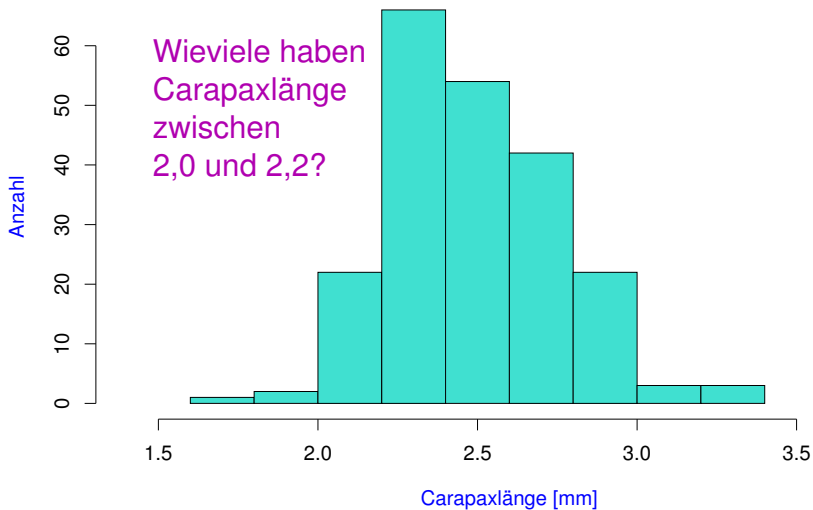


Eine Möglichkeit der graphischen  
Darstellung:  
das Histogramm

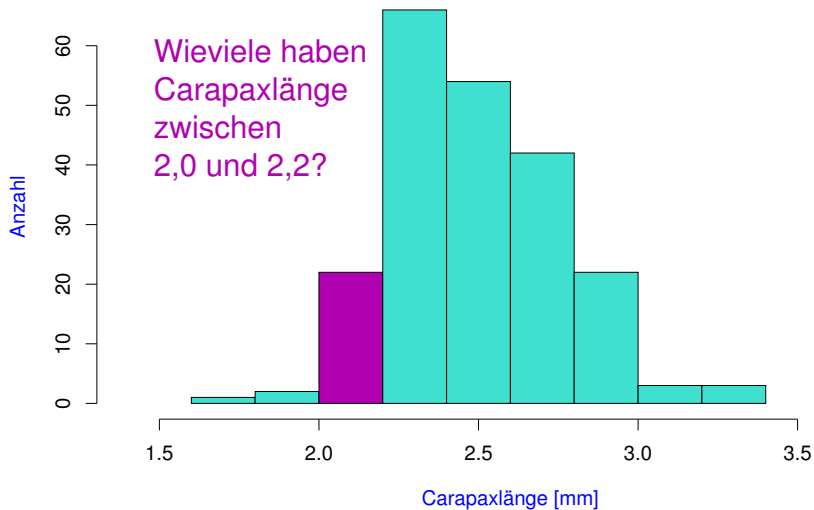
## Nichteiertragende Weibchen am 6. Sept. '88, n=215



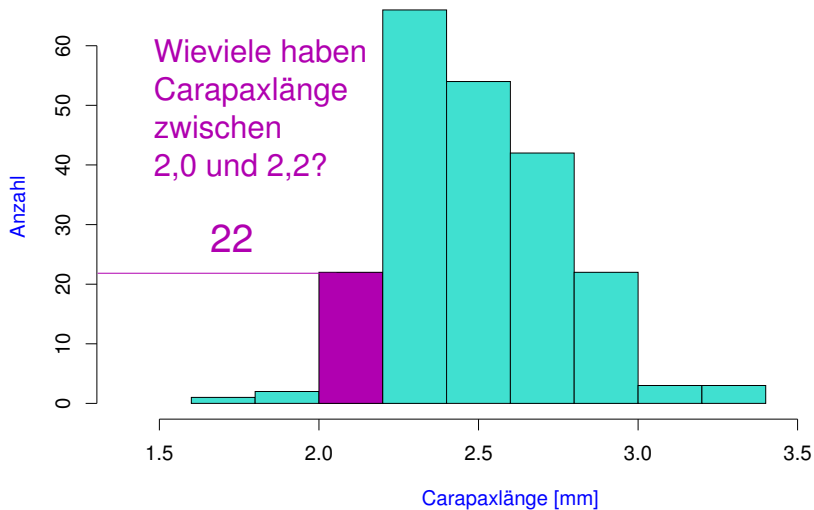
## Nichteiertragende Weibchen am 6. Sept. '88, n=215



## Nichteiertragende Weibchen am 6. Sept. '88, n=215

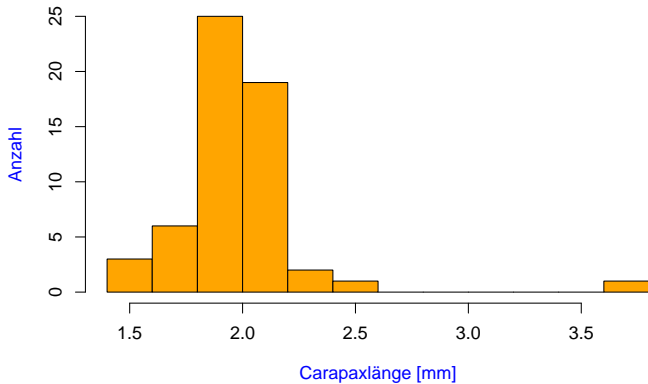


## Nichteiertragende Weibchen am 6. Sept. '88, n=215



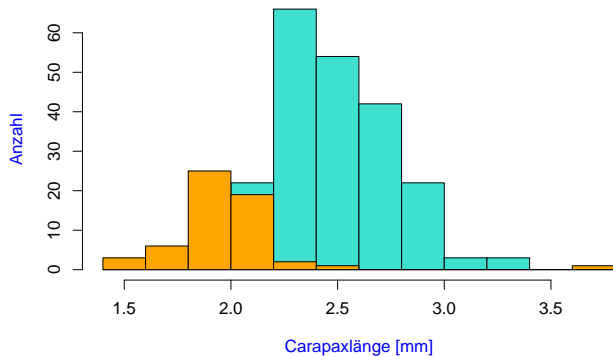
# Analoge Daten zwei Monate später (3.11.88):

### Nichteiertragende Weibchen am 3. Nov. '88, n=57



# Vergleich der beiden Verteilungen

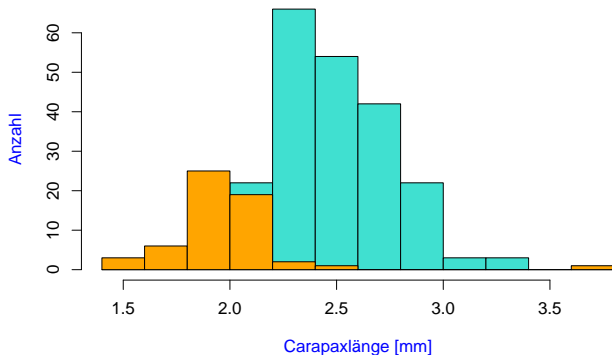
## Nichteiertragende Weibchen





# Vergleich der beiden Verteilungen

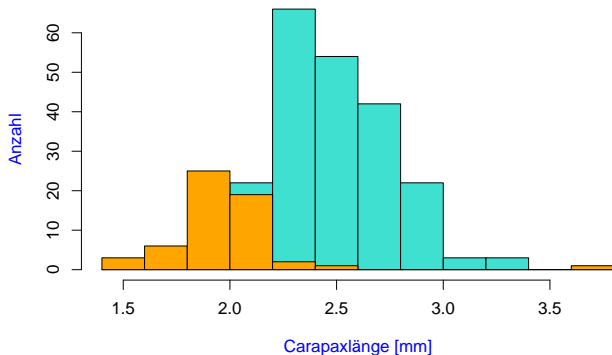
## Nichteiertragende Weibchen



Problem: ungleiche Stichprobenumfänge: 6.Sept:  $n = 215$   
3.Nov :  $n = 57$

# Vergleich der beiden Verteilungen

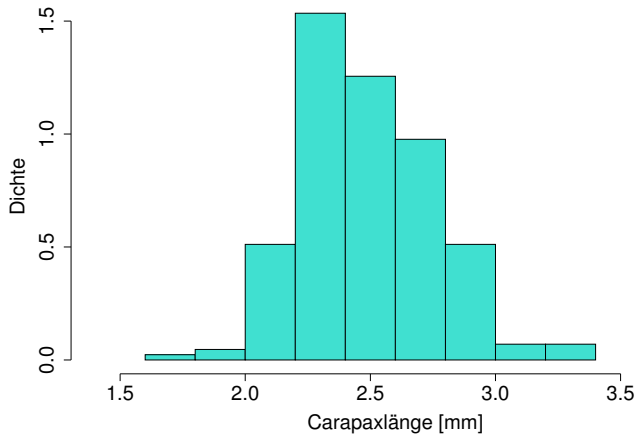
## Nichteiertragende Weibchen



Problem: ungleiche Stichprobenumfänge: 6.Sept:  $n = 215$   
 3.Nov :  $n = 57$

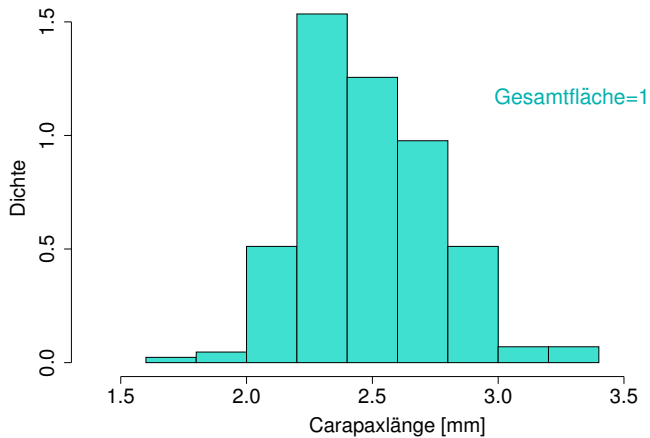
Idee: stauche vertikale Achse so, dass Gesamtfläche = 1.

Nichteiertragende Weibchen am 6. Sept. '88, n=215

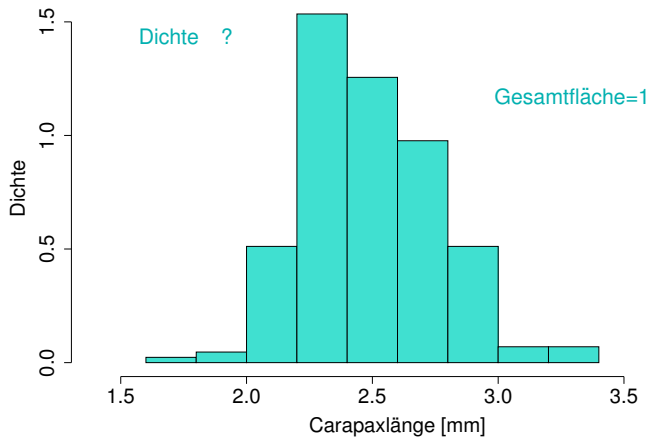


Die neue  
vertikale Koordinate  
ist jetzt eine  
**Dichte**  
(engl. **density**).

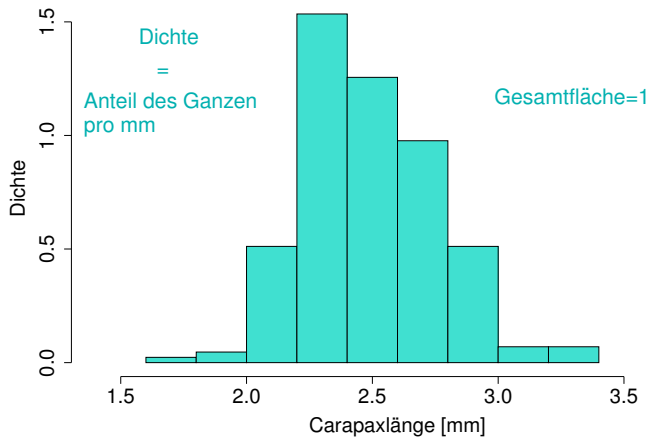
Nichteiertragende Weibchen am 6. Sept. '88, n=215



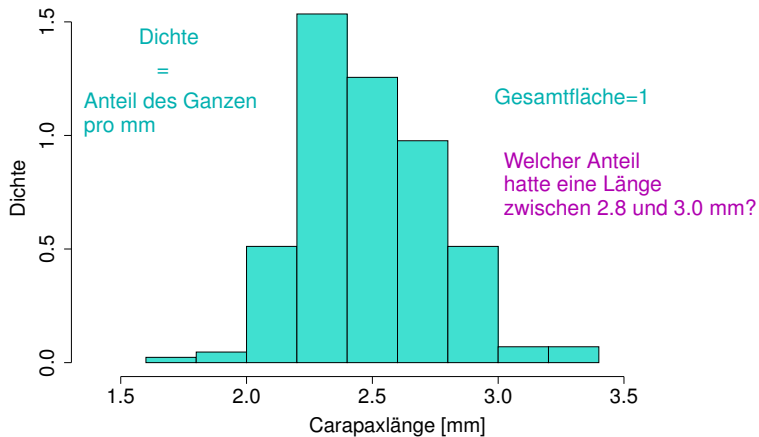
Nichteiertragende Weibchen am 6. Sept. '88, n=215



Nichteiertragende Weibchen am 6. Sept. '88, n=215

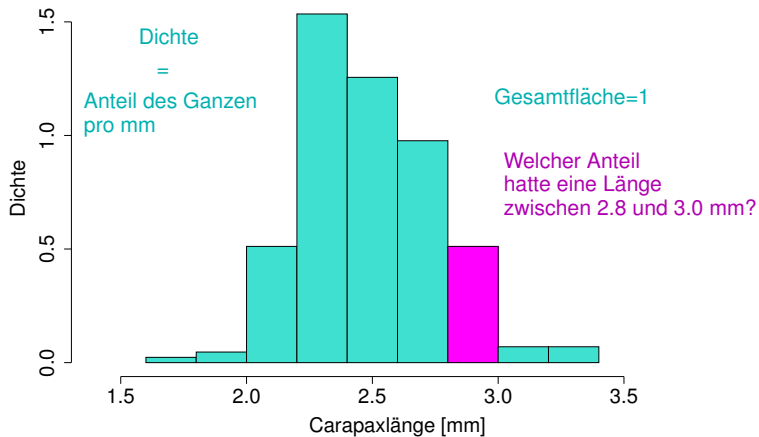


Nichteiertragende Weibchen am 6. Sept. '88, n=215

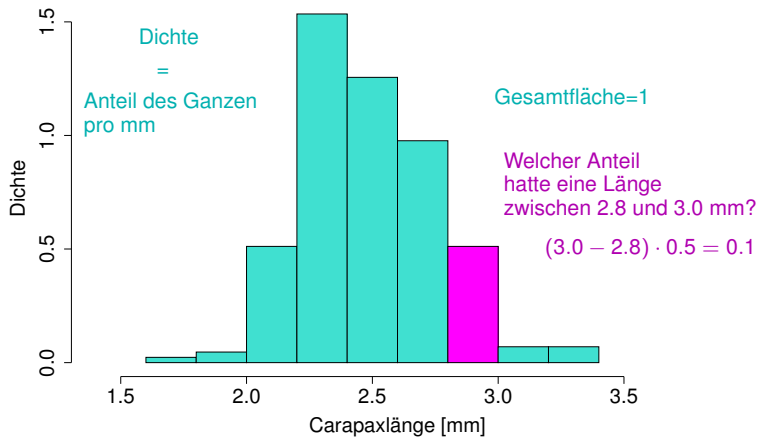




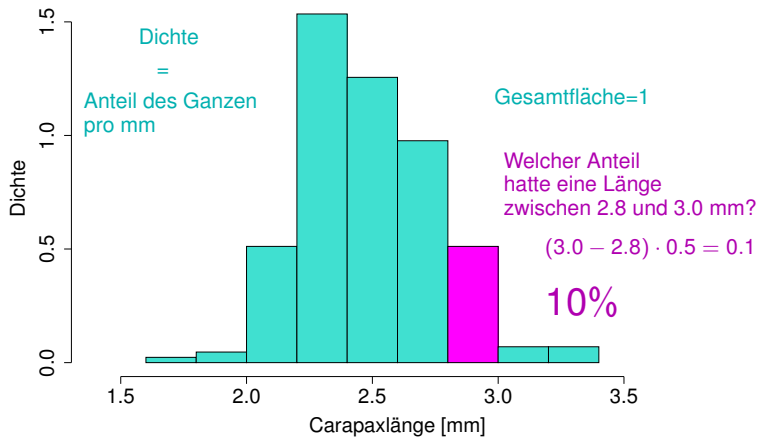
Nichteiertragende Weibchen am 6. Sept. '88, n=215



Nichteiertragende Weibchen am 6. Sept. '88, n=215



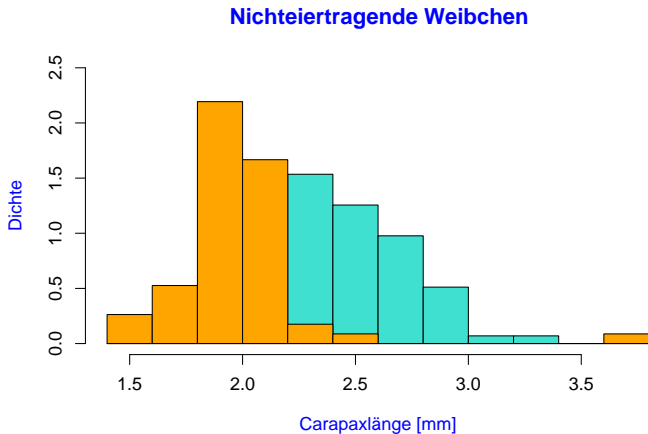
Nichteiertragende Weibchen am 6. Sept. '88, n=215



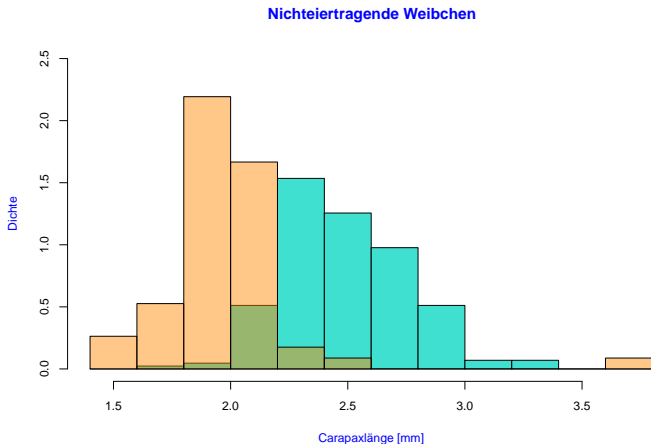
Die beiden Histogramme sind jetzt  
vergleichbar

Die beiden Histogramme sind jetzt  
vergleichbar  
(sie haben dieselbe Gesamtfläche).

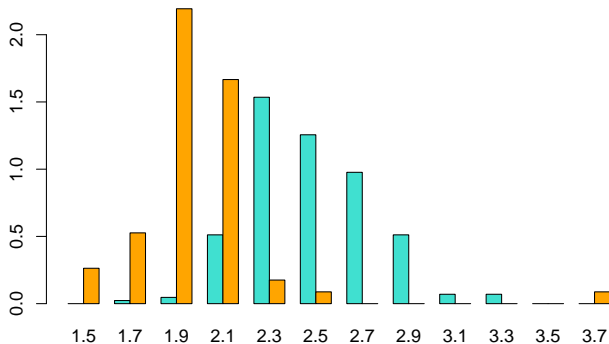
# Versuche, die Histogramme zusammen zu zeigen:



# Versuche, die Histogramme zusammen zu zeigen:

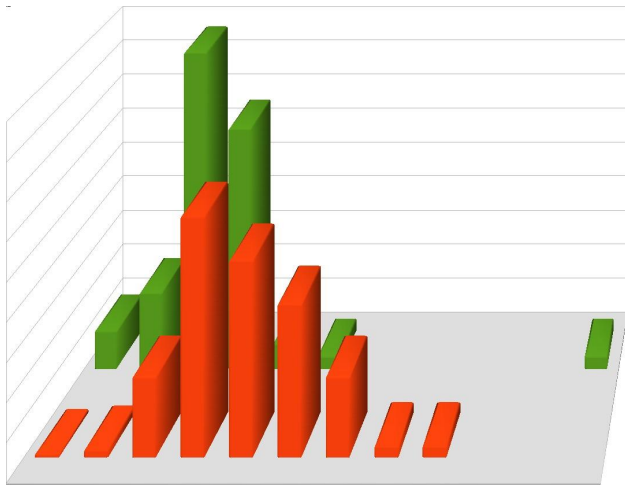


# Versuche, die Histogramme zusammen zu zeigen:





# Versuche, die Histogramme zusammen zu zeigen:



# Vorschlag

Total abgefahrene 3D-Plots können in der Werbung nützlich sein

# Vorschlag

Total abgefahrene 3D-Plots können in der Werbung nützlich sein,  
für die Wissenschaft sind einfache und klare 2D-Darstellungen  
meistens angemessener.

# Problem

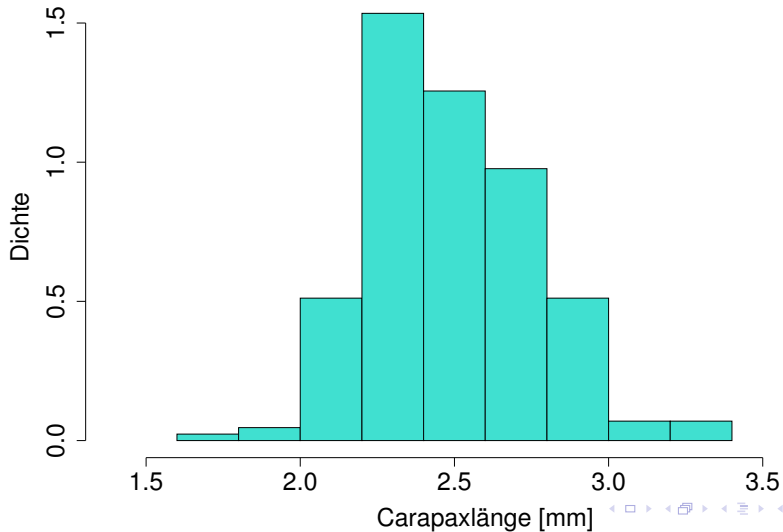
Histogramme kann man nicht ohne weiteres  
in demselben Graphen  
darstellen,

# Problem

Histogramme kann man nicht ohne weiteres  
in demselben Graphen  
darstellen,  
weil sie einander  
überdecken würden.

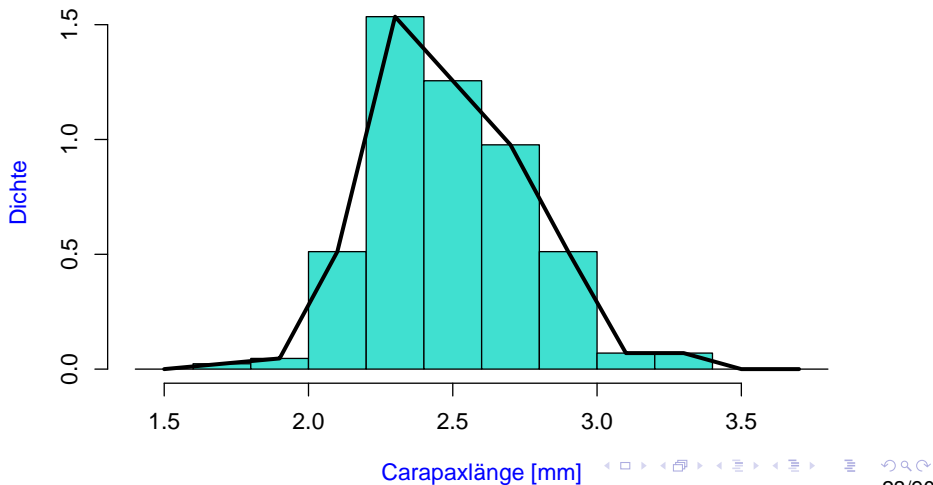
# Einfache und klare Lösung: Dichtepolygone

Nichteiertragende Weibchen am 6. Sept. '88, n=215



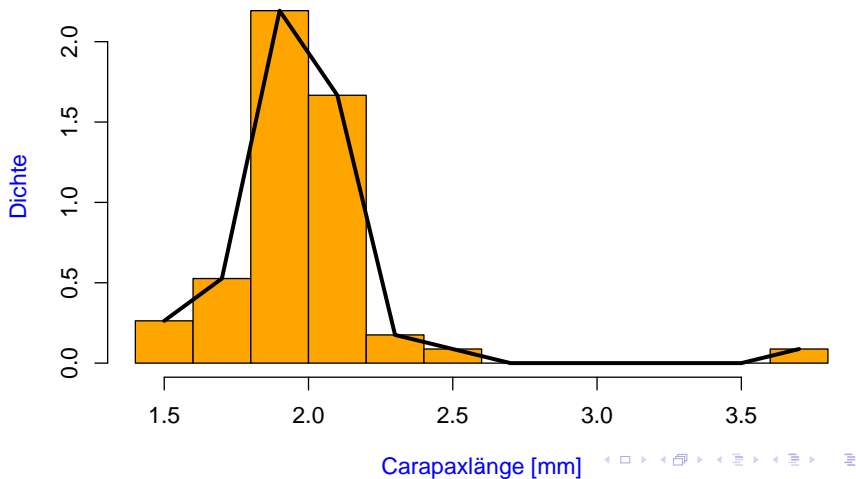
# Einfache und klare Lösung: Dichtepolygone

Nichteiertragende Weibchen am 6. Sept. '88, n=215



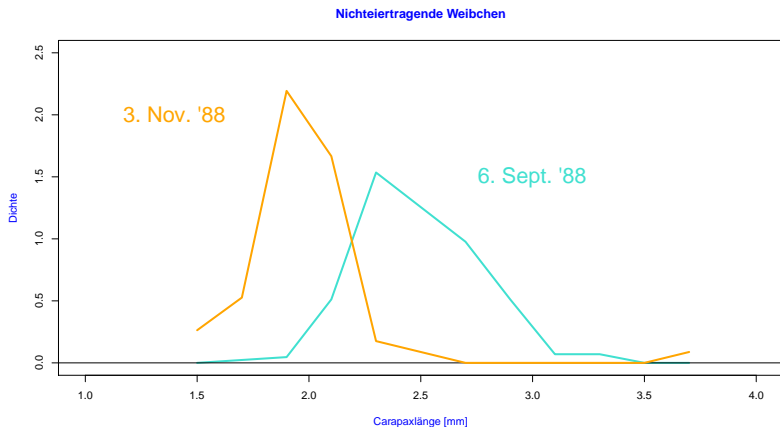
# Einfache und klare Lösung: Dichtepolygone

Nichteiertragende Weibchen am 3. Nov. '88, n=57

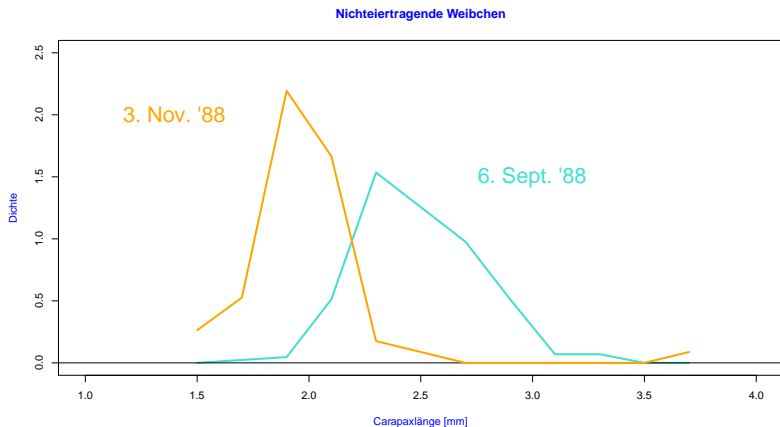




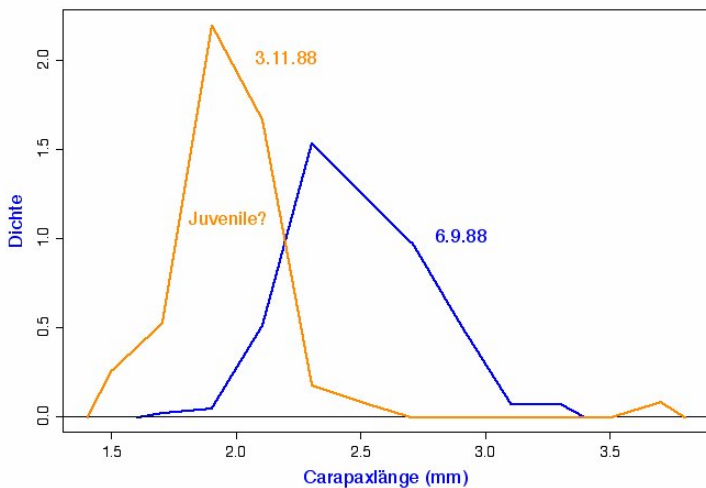
# Zwei und mehr Dichtepolygone in einem Plot

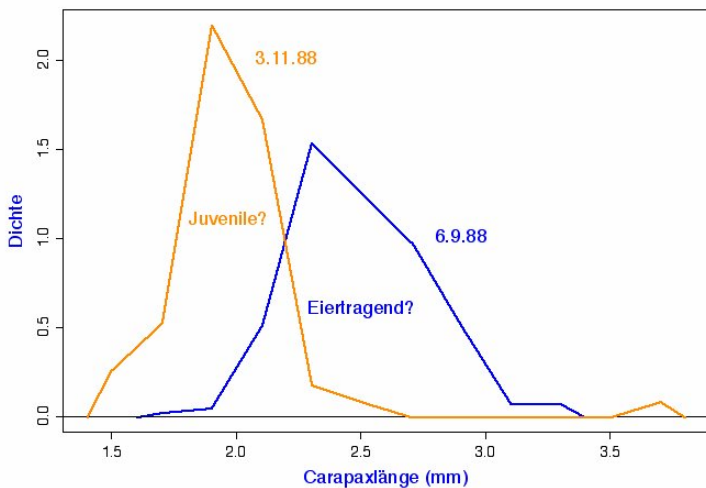


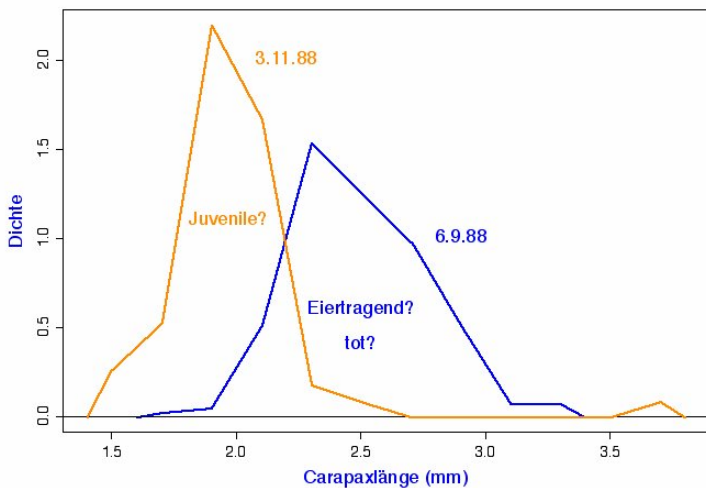
# Zwei und mehr Dichtepolygone in einem Plot



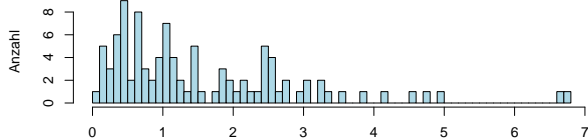
Biologische Interpretation der Verschiebung?

*Nichteiertragende Weibchen 6.9.88 und 3.11.88*

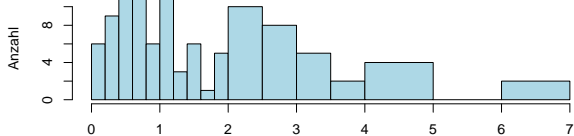
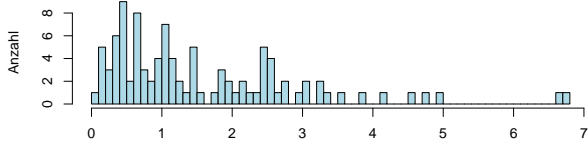
*Nichteiertragende Weibchen 6.9.88 und 3.11.88*

*Nichteiertragende Weibchen 6.9.88 und 3.11.88*

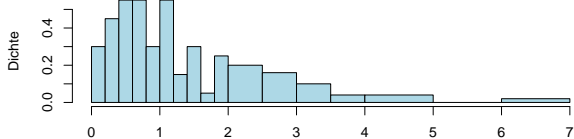
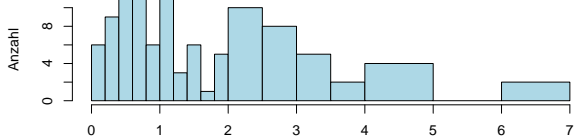
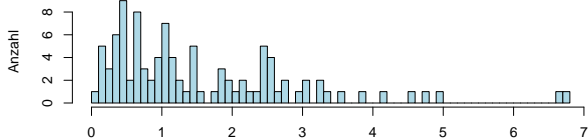
# Anzahl vs. Dichte



# Anzahl vs. Dichte

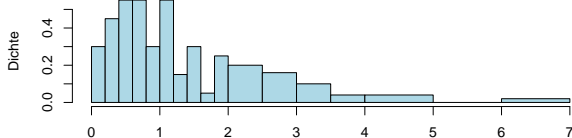
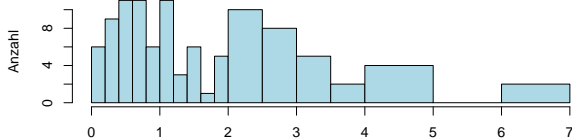
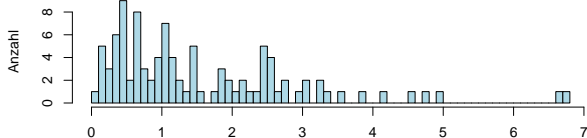


# Anzahl vs. Dichte





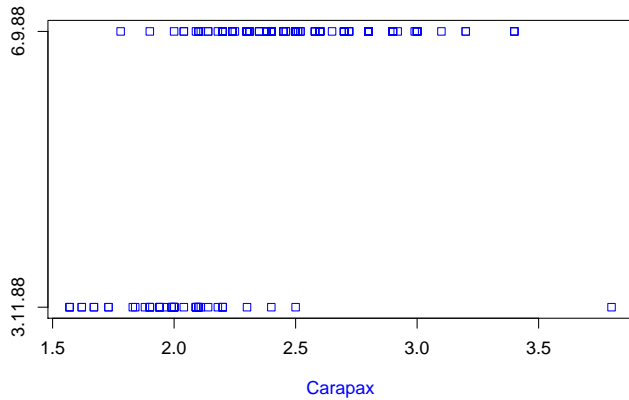
# Anzahl vs. Dichte

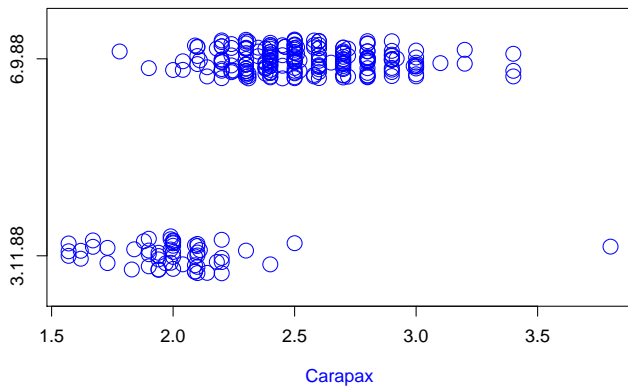


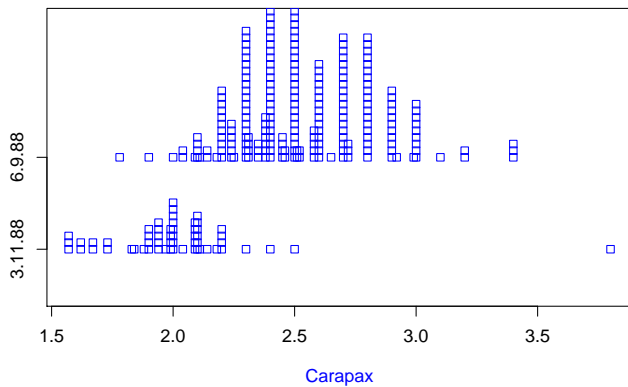
Also: Bei Histogrammen mit ungleichmäßiger Unterteilung immer Dichten verwenden!

# Inhalt

- 1 Ansatz der Statistik
- 2 **Graphische Darstellungen**
  - Histogramme und Dichtepolygone
  - **Stripcharts**
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R







Histogramme und Dichtepolygone  
geben  
ein ausführliches Bild  
eines Datensatzes.

Histogramme und Dichtepolygone  
geben  
ein ausführliches Bild  
eines Datensatzes.

Manchmal zu ausführlich.

# Inhalt

- 1 Ansatz der Statistik
- 2 **Graphische Darstellungen**
  - Histogramme und Dichtepolygone
  - Stripcharts
  - **Boxplots**
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R



# Zu viel Information erschwert den Überblick



Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum

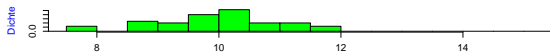
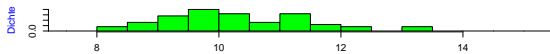
# Zu viel Information erschwert den Überblick

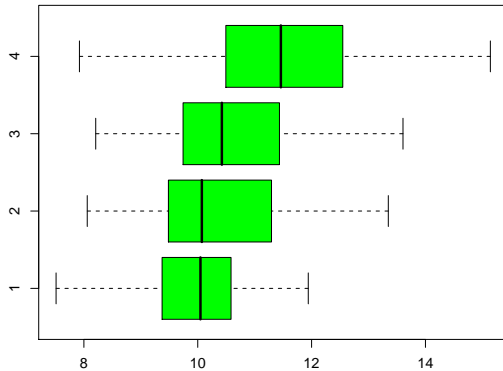


Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum Baum

Wald?

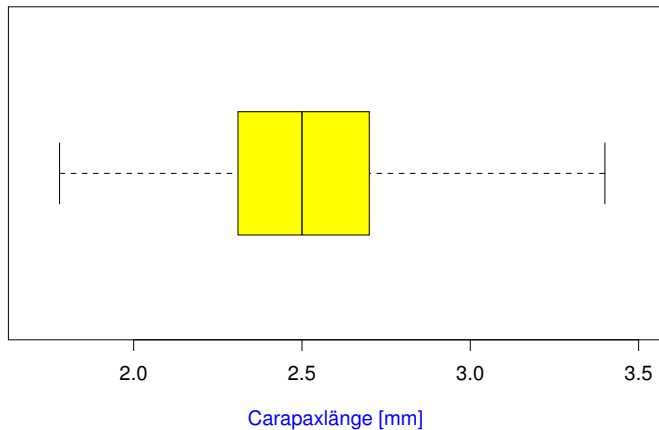
# Beispiel: Vergleich von mehreren Gruppen





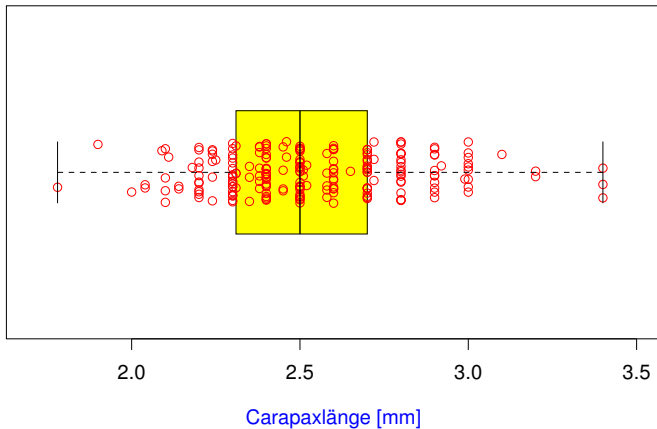
# Der Boxplot

## Boxplot, einfache Ausführung



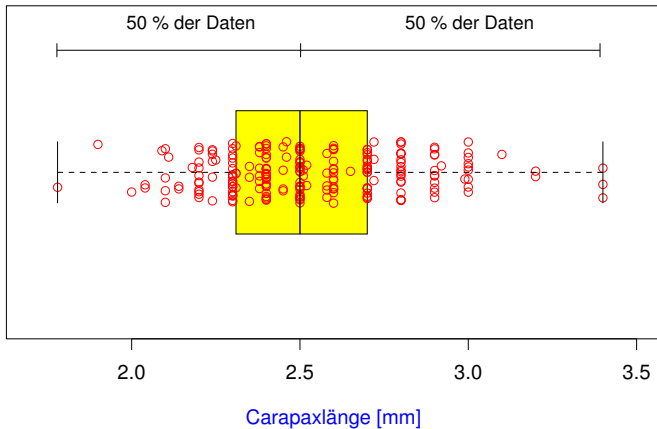
# Der Boxplot

## Boxplot, einfache Ausführung



# Der Boxplot

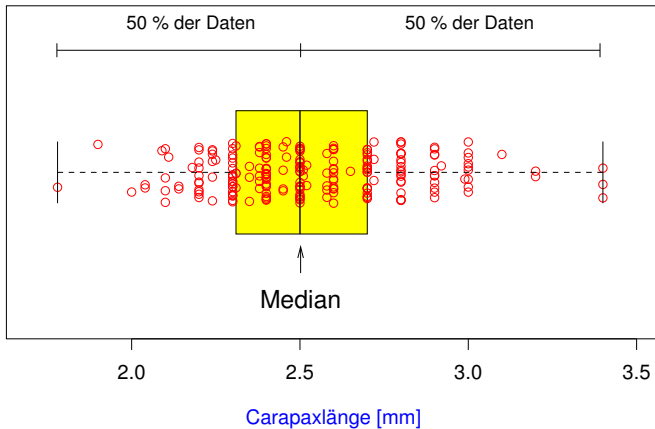
## Boxplot, einfache Ausführung





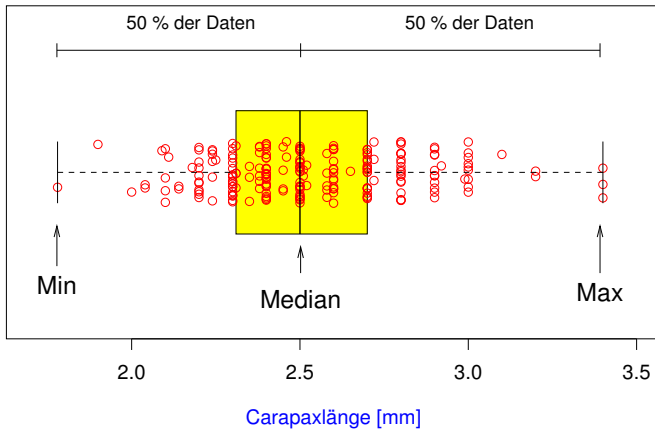
# Der Boxplot

## Boxplot, einfache Ausführung



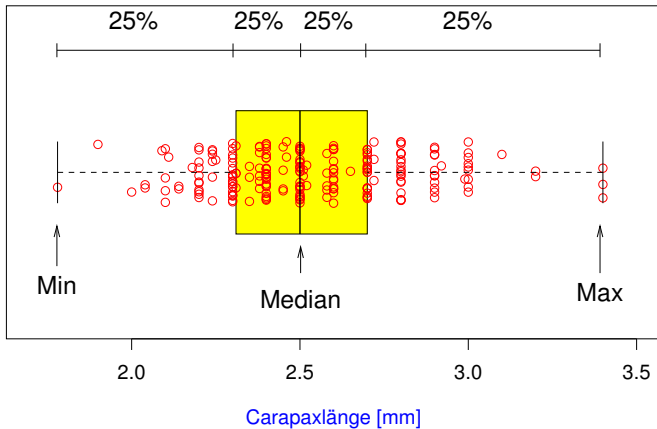
# Der Boxplot

## Boxplot, einfache Ausführung



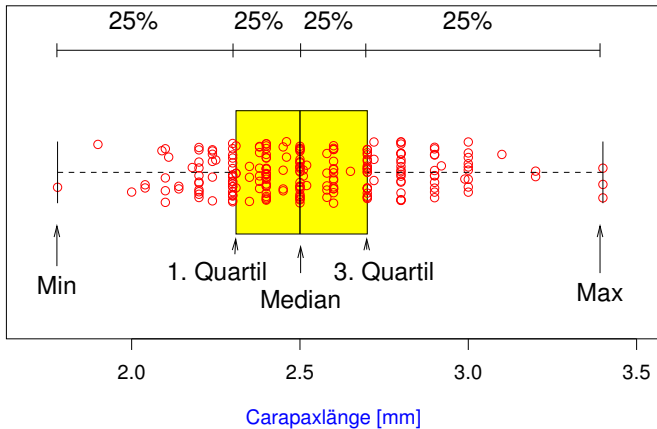
# Der Boxplot

## Boxplot, einfache Ausführung



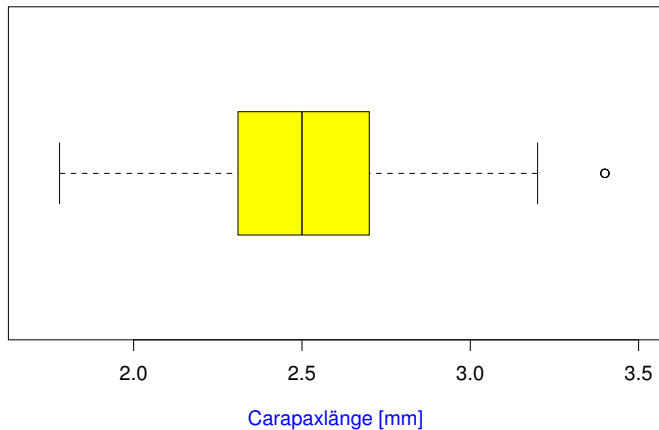
# Der Boxplot

## Boxplot, einfache Ausführung



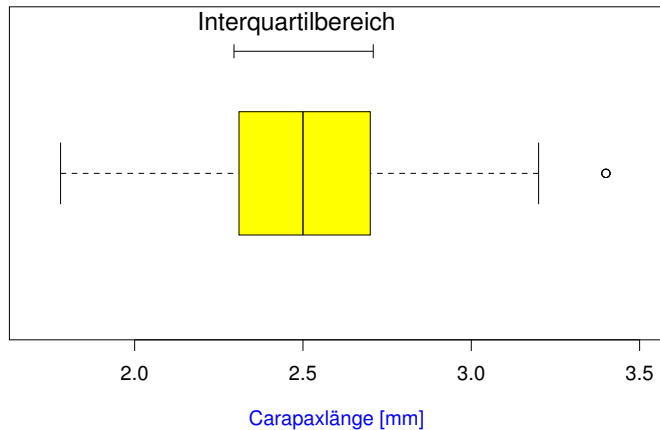
# Der Boxplot

## Boxplot, Standardausführung



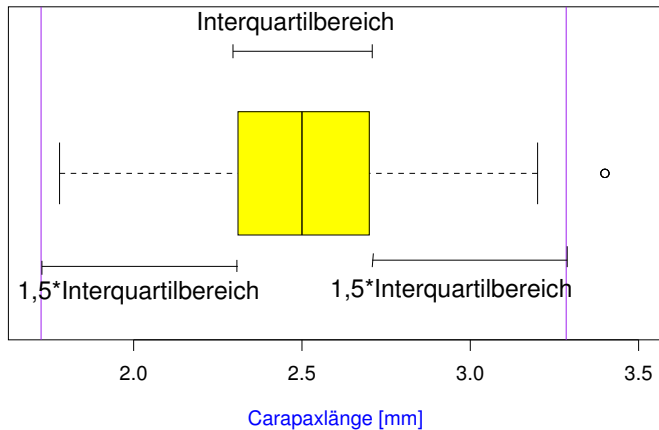
# Der Boxplot

## Boxplot, Standardausführung



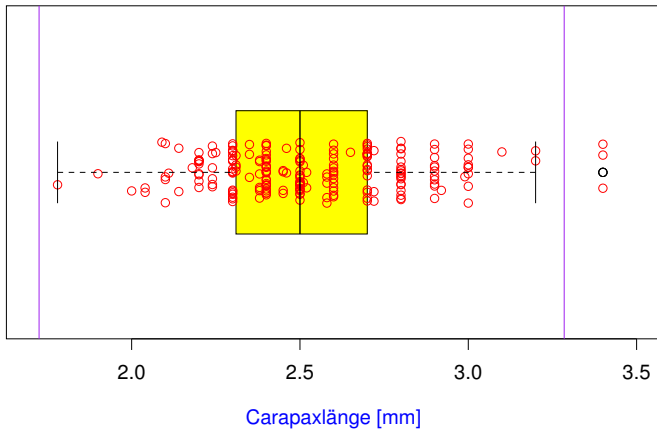
# Der Boxplot

## Boxplot, Standardausführung



# Der Boxplot

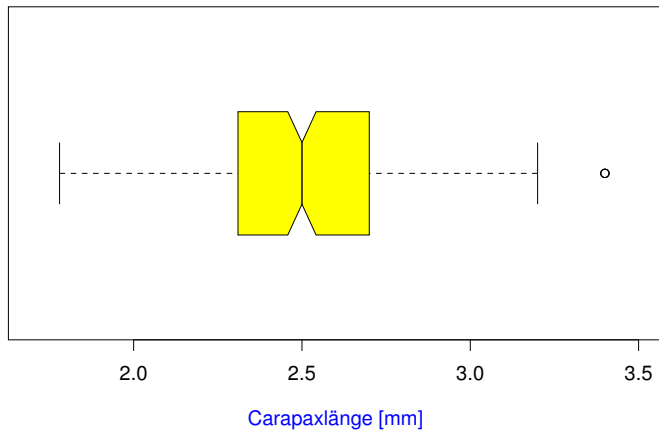
## Boxplot, Standardausführung





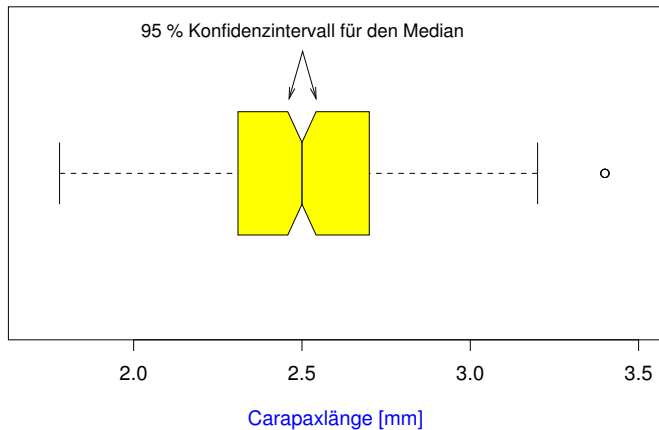
# Der Boxplot

## Boxplot, Profiausstattung



# Der Boxplot

## Boxplot, Profiausstattung



# Fazit

- 1 Histogramme erlauben einen detaillierten Blick auf die Daten

# Fazit

- 1 Histogramme erlauben einen detaillierten Blick auf die Daten
- 2 Dichtepolygone erlauben Vergleiche zwischen vielen Verteilungen

# Fazit

- 1 Histogramme erlauben einen detaillierten Blick auf die Daten
- 2 Dichtepolygone erlauben Vergleiche zwischen vielen Verteilungen
- 3 Boxplot können große Datenmengen vereinfacht zusammenfassen

# Fazit

- 1 Histogramme erlauben einen detaillierten Blick auf die Daten
- 2 Dichtepolygone erlauben Vergleiche zwischen vielen Verteilungen
- 3 Boxplot können große Datenmengen vereinfacht zusammenfassen
- 4 Bei kleinen Datenmengen eher Stripcharts angemessen

# Fazit

- 1 Histogramme erlauben einen detaillierten Blick auf die Daten
- 2 Dichtepolygone erlauben Vergleiche zwischen vielen Verteilungen
- 3 Boxplot können große Datenmengen vereinfacht zusammenfassen
- 4 Bei kleinen Datenmengen eher Stripcharts angemessen
- 5 Vorsicht mit Tricks wie 3D oder halbtransparenten Farben

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen**
  - Median und andere Quartile, (empirische) Quantile**
  - Mittelwert und Standardabweichung**
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R



Es ist oft möglich,  
das Wesentliche  
an einer Stichprobe  
  
mit ein paar Zahlen  
zusammenzufassen.

Wesentlich:

1. Wie groß?

2. Wie variabel?

Wesentlich:

1. Wie groß?

Lageparameter

2. Wie variabel?

Wesentlich:

1. Wie groß?

Lageparameter

2. Wie variabel?

Streuungsparameter

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen**
  - Median und andere Quartile, (empirische) Quantile**
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

Eine Möglichkeit  
kennen wir schon  
aus dem Boxplot:

# Lageparameter

## Der Median

# Lageparameter

## Der Median

# Streuungsparameter



# Lageparameter

## Der Median

# Streuungsparameter

## Der Quartilabstand ( $Q_3 - Q_1$ )

## Der **Median**<sup>1</sup>:

die Hälfte der Beobachtungen sind kleiner,  
die Hälfte sind größer.

---

<sup>1</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition)

Der **Median**<sup>1</sup>:

die Hälfte der Beobachtungen sind kleiner,  
die Hälfte sind größer.

Der Median ist  
das **50%-Quantil**  
der Daten.

---

<sup>1</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition)

# Die Quartile

Das erste Quartil<sup>2</sup>,  $Q_1$ :

---

<sup>2</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition)

# Die Quartile

Das erste Quartil<sup>2</sup>,  $Q_1$ :  
ein Viertel der Beobachtungen  
sind kleiner,  
drei Viertel sind größer.

---

<sup>2</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition) > < ≡ > ≡

# Die Quartile

Das erste Quartil<sup>2</sup>,  $Q_1$ :  
ein Viertel der Beobachtungen  
sind kleiner,  
drei Viertel sind größer.

$Q_1$  ist das  
**25%-Quantil**  
der Daten.

---

<sup>2</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition) > < ≡ ≡ ≡

# Die Quartile

Das dritte Quartil<sup>3</sup>,  $Q_3$ :

---

<sup>3</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition) > < ≡ ≡ ≡

# Die Quartile

Das dritte Quartil<sup>3</sup>,  $Q_3$ :  
drei Viertel der Beobachtungen  
sind kleiner,  
ein Viertel sind größer.

---

<sup>3</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition) > < ≡ ≡ ≡



# Die Quartile

Das dritte Quartil<sup>3</sup>,  $Q_3$ :  
drei Viertel der Beobachtungen  
sind kleiner,  
ein Viertel sind größer.

$Q_3$  ist das  
75%-Quantil  
der Daten.

---

<sup>3</sup>„saloppe“ Definition (wir sehen gleich die präzise Definition) ▶ ◀ ≡ ≡ ≡

# (Empirische) Quantile, allgemein

Seien  $n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$  gegeben,  
 $\alpha \in (0, 1)$ .

$q$  ist (ein)  $\alpha$ -Quantil der  $n$  Beobachtungswerte, wenn gilt

$$\frac{1}{n} |\{1 \leq i \leq n : x_i \leq q\}| \geq \alpha \text{ und } \frac{1}{n} |\{1 \leq i \leq n : x_i \geq q\}| \geq 1 - \alpha.$$

# (Empirische) Quantile, allgemein

Seien  $n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$  gegeben,  
 $\alpha \in (0, 1)$ .

$q$  ist (ein)  $\alpha$ -Quantil der  $n$  Beobachtungswerte, wenn gilt

$$\frac{1}{n} |\{1 \leq i \leq n : x_i \leq q\}| \geq \alpha \text{ und } \frac{1}{n} |\{1 \leq i \leq n : x_i \geq q\}| \geq 1 - \alpha.$$

**Bem.:** Im Allgemeinen ist ein  $\alpha$ -Quantil nicht eindeutig:

Seien  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$  die der Größe nach sortierten  
Werte.

Wenn  $\alpha = \frac{k}{n}$  mit  $1 \leq k < n$ , so ist jeder Wert  $q \in [x_{(k)}, x_{(k+1)}]$  ein  
 $\alpha$ -Quantil,

denn  $|\{i : x_i \leq x_{(k)}\}| \geq k$ ,  $|\{i : x_i \geq x_{(k)}\}| \geq n - k + 1$ .

Wenn  $n\alpha \notin \{1, \dots, n-1\}$ , so ist das  $\alpha$ -Quantil der Wert  $x_{(k)}$  mit  
 $k = \lceil \alpha n \rceil$ .

# (Empirische) Quantile, allgemein II

$n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$  gegeben,  $\alpha \in (0, 1)$ .

(ein)  $\alpha$ -Quantil  $q$  der  $n$  Beobachtungswerte erfüllt

$$\frac{1}{n} |\{1 \leq i \leq n : x_i \leq q\}| \geq \alpha \text{ und } \frac{1}{n} |\{1 \leq i \leq n : x_i \geq q\}| \geq 1 - \alpha.$$

# (Empirische) Quantile, allgemein II

$n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$  gegeben,  $\alpha \in (0, 1)$ .  
(ein)  $\alpha$ -Quantil  $q$  der  $n$  Beobachtungswerte erfüllt

$$\frac{1}{n} |\{1 \leq i \leq n : x_i \leq q\}| \geq \alpha \text{ und } \frac{1}{n} |\{1 \leq i \leq n : x_i \geq q\}| \geq 1 - \alpha.$$

## Bem.:

- Die Definition passt zu unserer früheren Definition für Verteilungen, wenn man die *empirische Verteilung*  $\frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  betrachtet.
- In der Literatur (und auch in Statistik-Software) sind verschiedene Interpolationen üblich, um „das“  $\alpha$ -Quantil stetig in  $\alpha$  zu machen.  
(In R siehe etwa `help(quantile)`, es sind 9 Varianten implementiert.)
- Die Uneindeutigkeit des  $\alpha$ -Quantils ist für halbwegs große  $n$  in der Praxis oft wenig von Belang.

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen**
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung**
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

$n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$

Am häufigsten werden benutzt:

Lageparameter

Der **Mittelwert**  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$

$n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$

Am häufigsten werden benutzt:

Lageparameter

Der Mittelwert  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$

Streuungsparameter

Die **Standardabweichung**  $s$  (bzw.  $\sigma$ )



$n$  (reelle) Beobachtungswerte  $x_1, x_2, \dots, x_n$

Am häufigsten werden benutzt:

Lageparameter

Der Mittelwert  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$

Streuungsparameter

Die Standardabweichung  $s$  (bzw.  $\sigma$ )

wobei

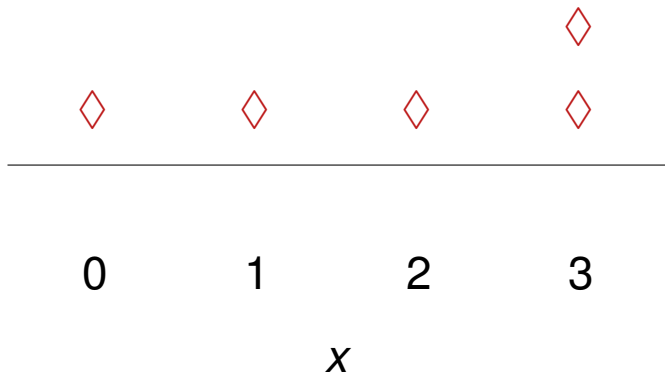
$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{die (empirische) Varianz}$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{die korrigierte Stichproben-Varianz}$$

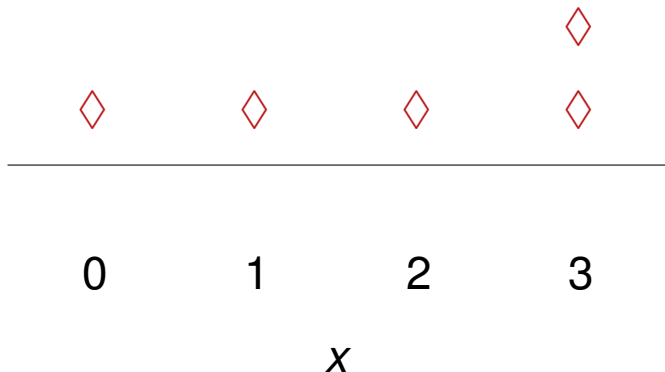
$$\left( = \frac{n}{n-1} \sigma^2 \right)$$

# Erinnerung: Geometrische Bedeutung des Mittelwerts Der Schwerpunkt

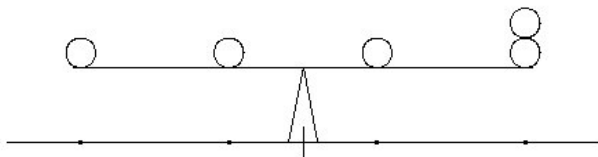
Wir stellen uns die Beobachtungen als  
gleich schwere Gewichte auf einer Waage  
vor:



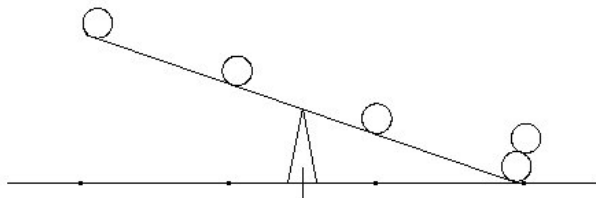
Wo muß der Drehpunkt sein, damit die Waage im Gleichgewicht ist?



$$m = 1,5 ?$$

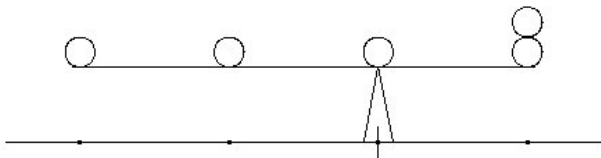


$$m = 1,5 ?$$

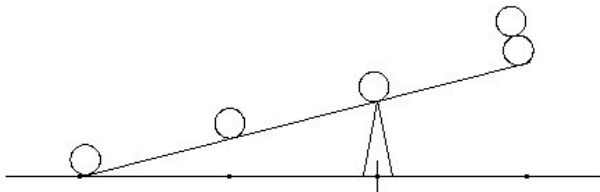


zu klein

$$m = 2 ?$$



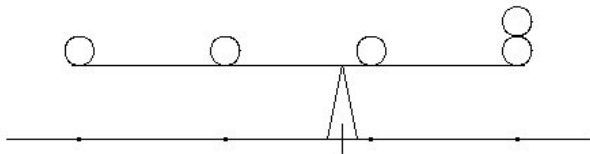
$$m = 2 ?$$



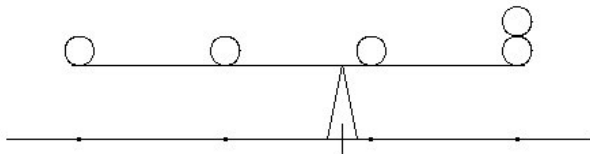
zu groß



$$m = 1,8 ?$$



$$m = 1,8 ?$$



richtig

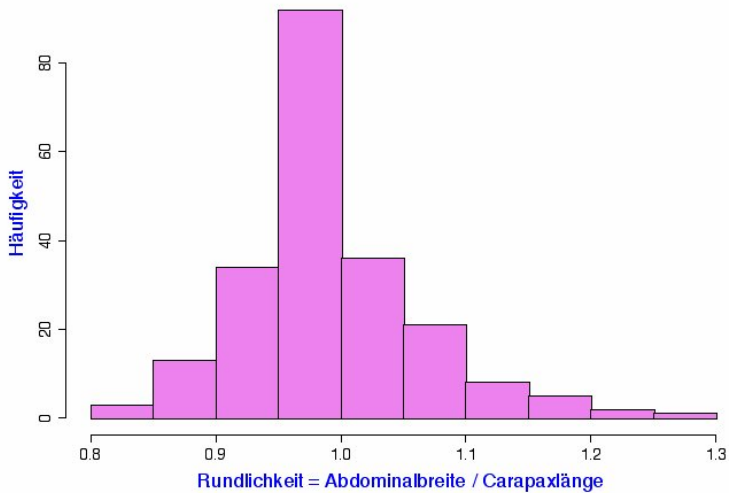
Oft kann man „mit dem bloßen Auge“ anhand eines Histogramms den Mittelwert gut einschätzen.

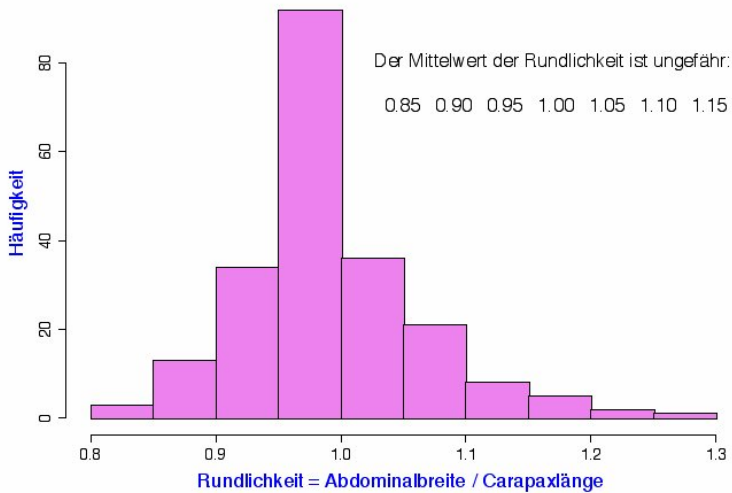
Beispiel: *Galathea intermedia*

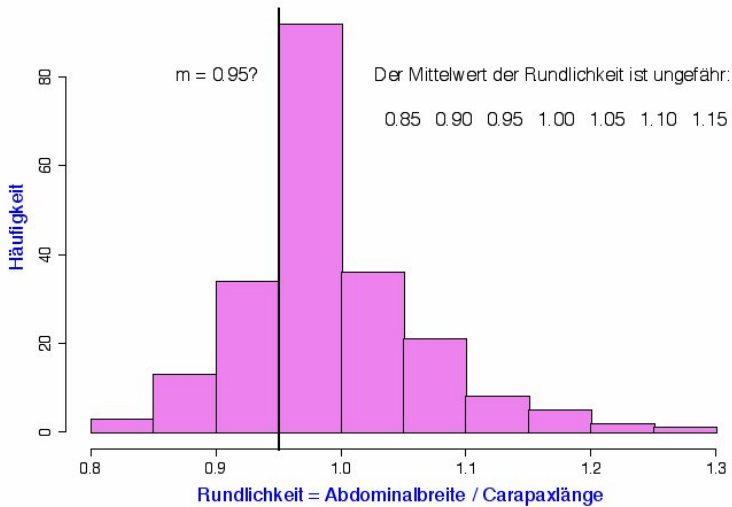
„Rundlichkeit“

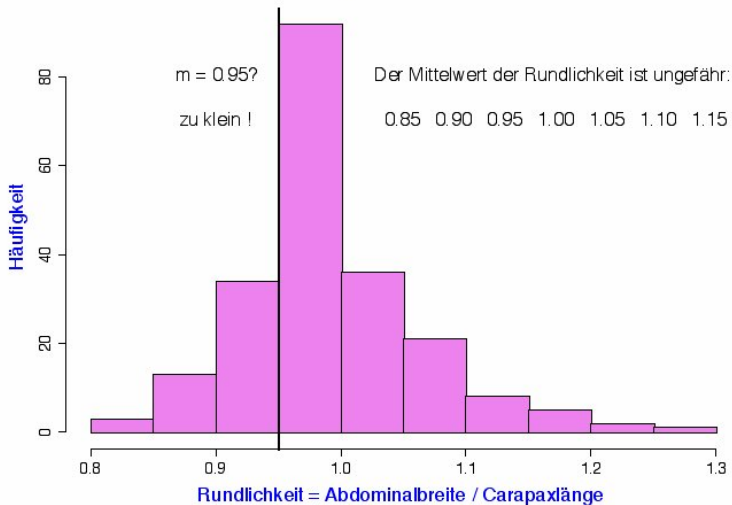
:=

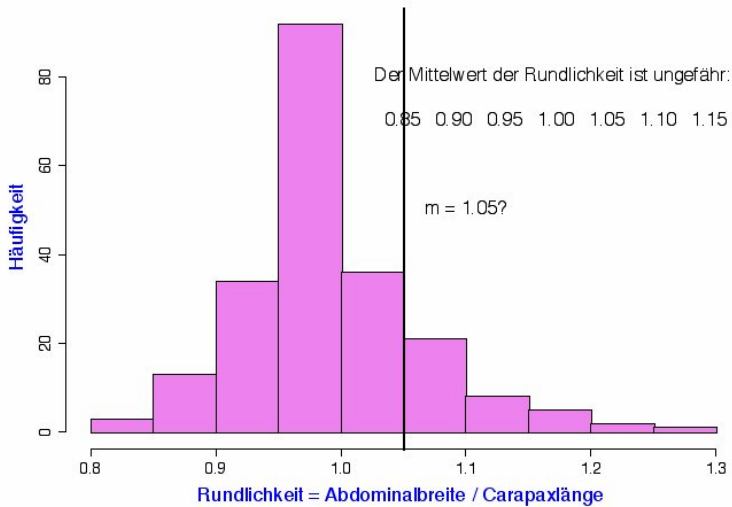
Abdominalbreite / Carapaxlänge

*Nichteiertragende Weibchen 6.9.88*

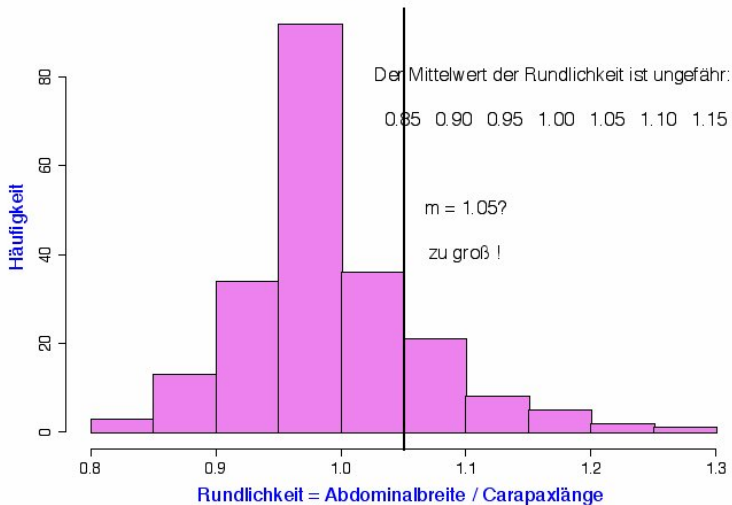
*Nichteiertragende Weibchen 6.9.88*

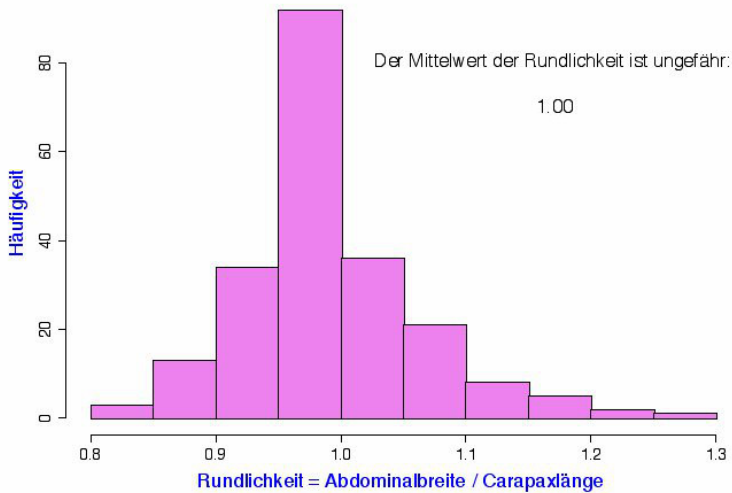
*Nichteiertragende Weibchen 6.9.88*

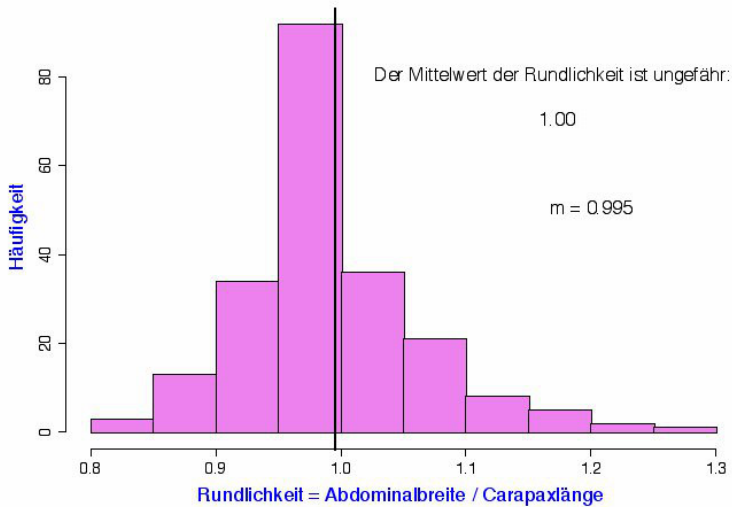
*Nichteiertragende Weibchen 6.9.88*

*Nichteiertragende Weibchen 6.9.88*

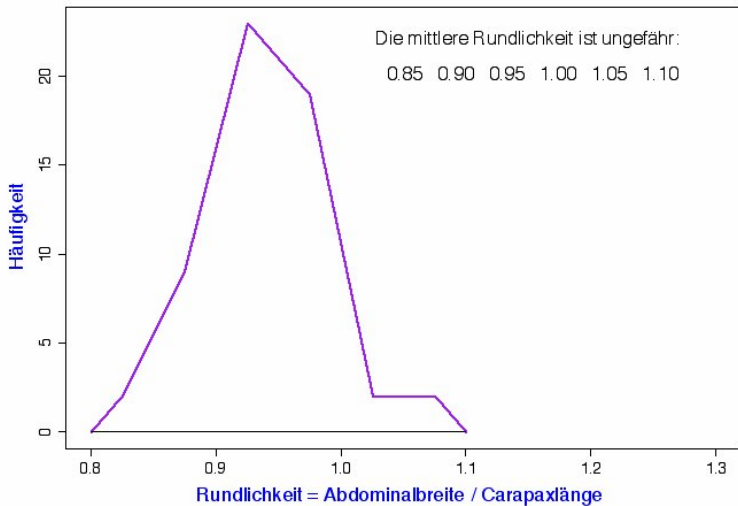


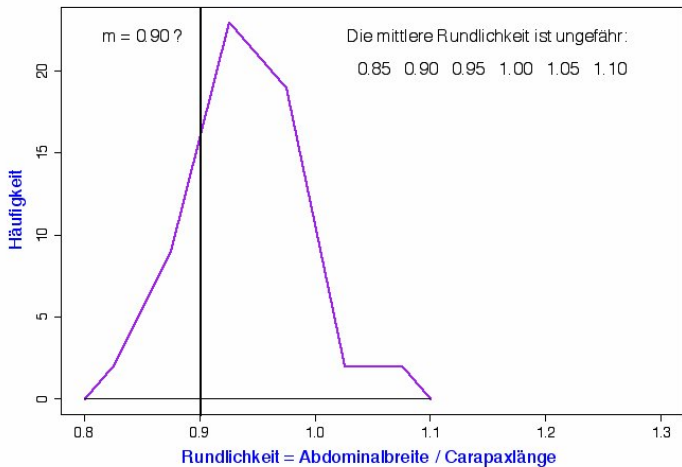
*Nichteiertragende Weibchen 6.9.88*

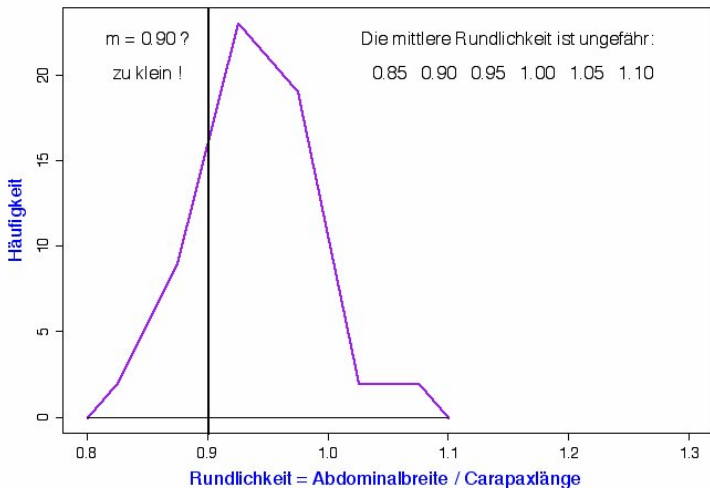
*Nichteiertragende Weibchen 6.9.88*

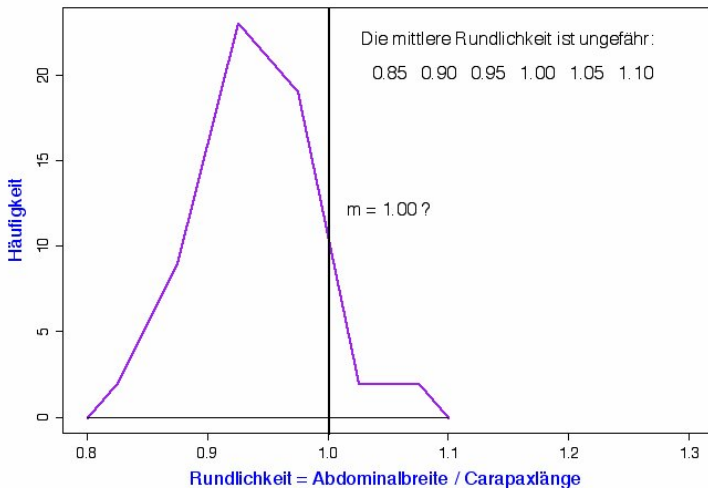
*Nichteiertragende Weibchen 6.9.88*

Beispiel:  
3.11.88

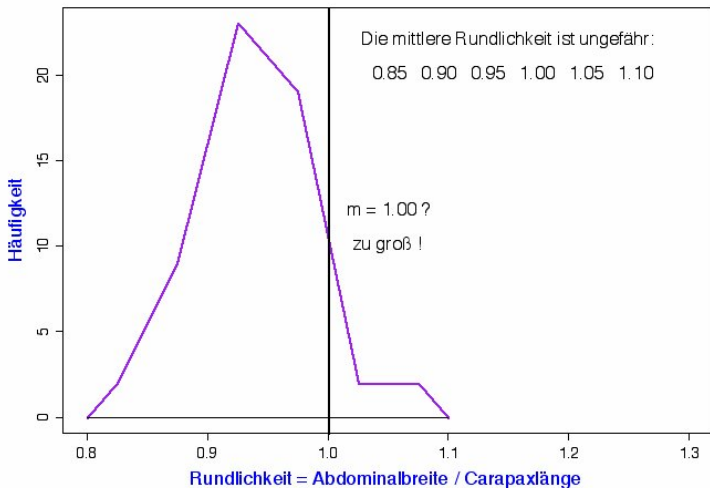
*Nichteiertragende Weibchen 3.11.88*

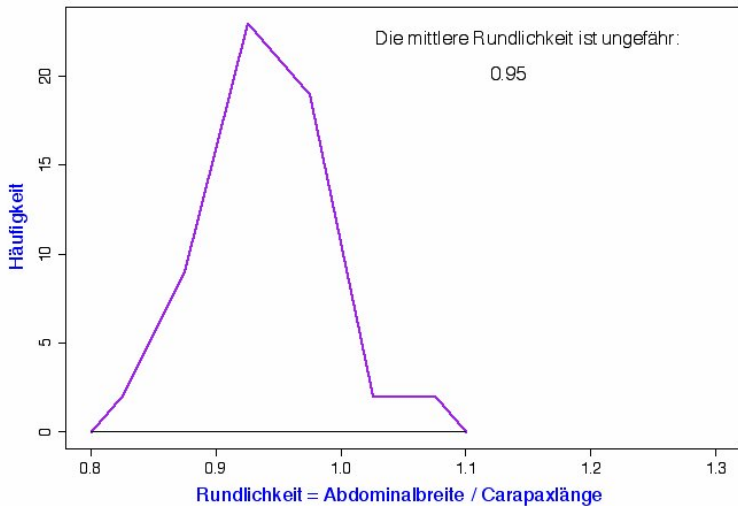
*Nichteiertragende Weibchen 3.11.88*

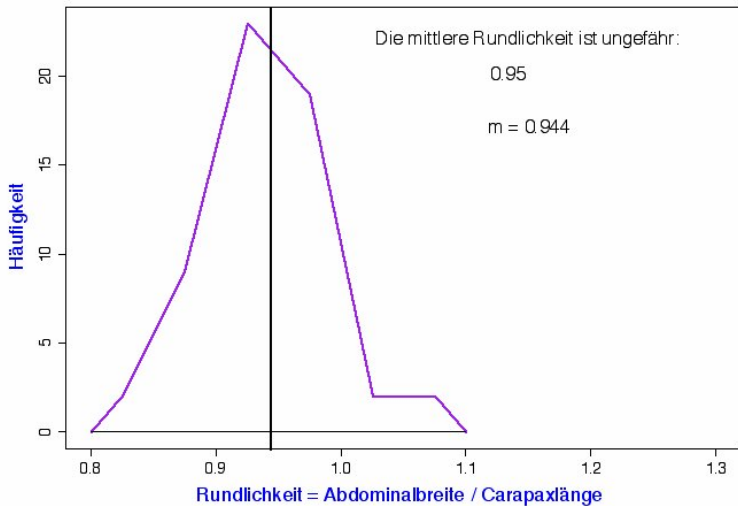
*Nichteiertragende Weibchen 3.11.88*

*Nichteiertragende Weibchen 3. 11. 88*



*Nichteiertragende Weibchen 3.11.88*

*Nichteiertragende Weibchen 3.11.88*

*Nichteiertragende Weibchen 3.11.88*

## Die Standardabweichung (auch: Streuung)

## Die Standardabweichung (auch: Streuung)

Wie weit weicht  
eine typische Beobachtung  
vom  
Mittelwert  
ab ?

# Mit $n$ oder $n - 1$ berechnen?

Die Standardabweichung  $\sigma$  eines Zufallsexperiments mit  $n$  gleichwahrscheinlichen Ausgängen  $x_1, \dots, x_n$  (z.B. Würfelwurf) ist definiert durch

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

# Mit $n$ oder $n - 1$ berechnen?

Die Standardabweichung  $\sigma$  eines Zufallsexperiments mit  $n$  gleichwahrscheinlichen Ausgängen  $x_1, \dots, x_n$  (z.B. Würfelwurf) ist definiert durch

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

Wenn es sich bei  $x_1, \dots, x_n$  um Beobachtungswerte in einer Stichprobe handelt, verwendet man eher

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

# s als Schätzer für $\sigma$

Wir werden sehen:

Wenn  $X_1, \dots, X_n$  u.i.v. *Zufallsvariablen* mit Varianz  $\text{Var}[X_1] = \sigma^2$ ,

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i,$$

so hat die *Zufallsvariable*

$$S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

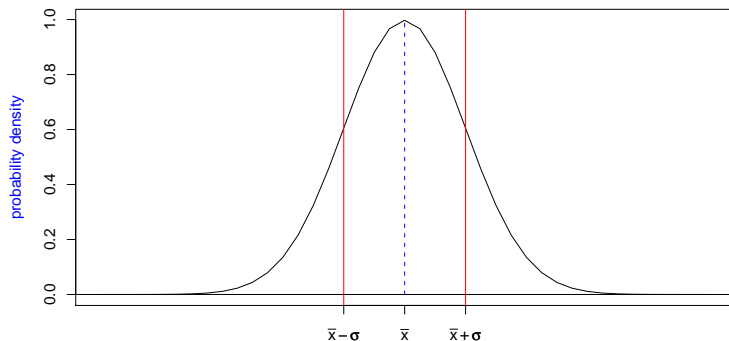
die Eigenschaft

$$\mathbb{E}[S^2] = \sigma^2.$$



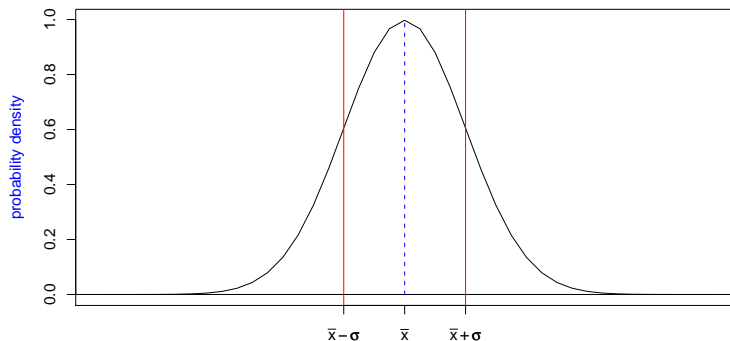
# Faustregel für die Standardabweichung

Bei ungefähr glockenförmigen (also eingipfligen und symmetrischen) Verteilungen liegen ca. 2/3 der Verteilung zwischen  $\bar{x} - \sigma$  und  $\bar{x} + \sigma$ .



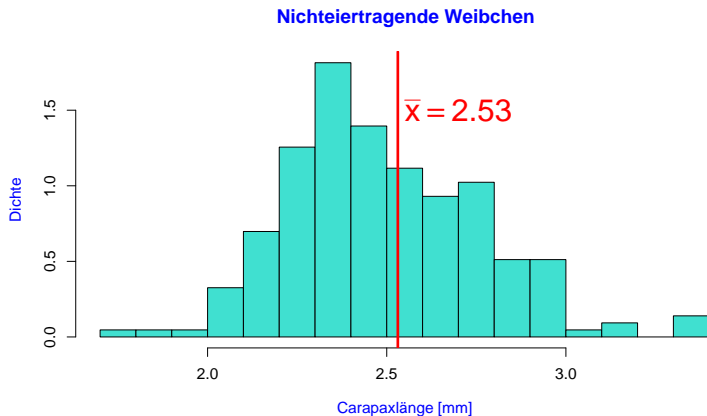
# Faustregel für die Standardabweichung

Bei ungefähr glockenförmigen (also eingipfligen und symmetrischen) Verteilungen liegen ca. 2/3 der Verteilung zwischen  $\bar{x} - \sigma$  und  $\bar{x} + \sigma$ .

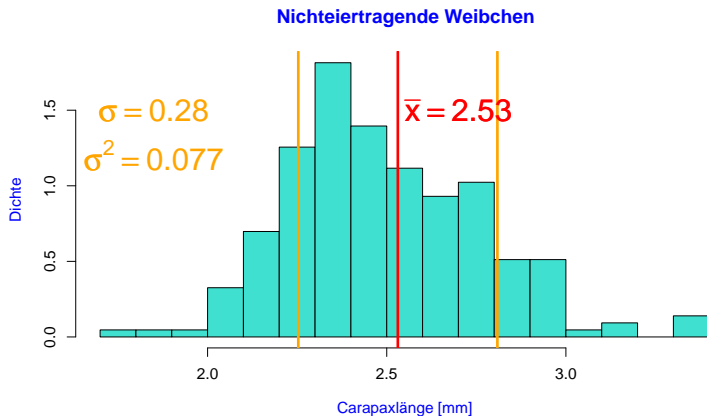


Oft kann man so die Standardabweichung „mit bloßem Auge“ abschätzen.

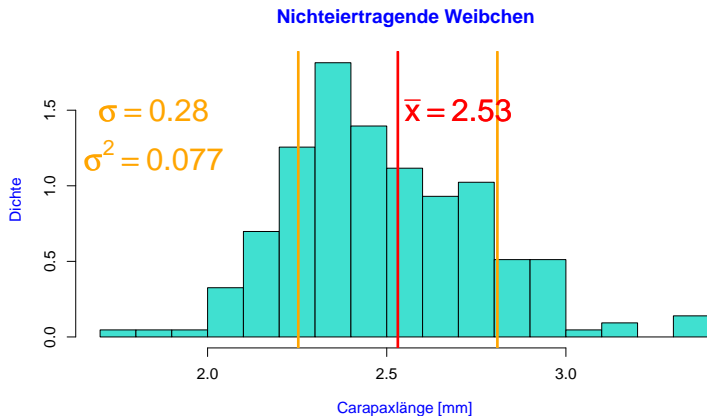
# Standardabweichung der Carapaxlängen nichteierttragender Weibchen vom 6.9.88



# Standardabweichung der Carapaxlängen nichteierttragender Weibchen vom 6.9.88



# Standardabweichung der Carapaxlängen nichteiertragender Weibchen vom 6.9.88



Hier liegt der Anteil zwischen  $\bar{x} - \sigma$  und  $\bar{x} + \sigma$  bei 72%.

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten**
  - **Beispiel: Wählerische Bachstelzen**
  - **Beispiel: Spiderman & Spiderwoman**
  - **Beispiel: Kupfertoleranz beim Roten Straußgras**
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

## Mittelwert und Standardabweichung. . .

- charakterisieren die Daten gut, falls deren Verteilung glockenförmig ist

## Mittelwert und Standardabweichung. . .


- charakterisieren die Daten gut, falls deren Verteilung glockenförmig ist
- und müssen andernfalls mit Vorsicht interpretiert werden.



## Mittelwert und Standardabweichung. . .

- charakterisieren die Daten gut, falls deren Verteilung glockenförmig ist
- und müssen andernfalls mit Vorsicht interpretiert werden.


Wir betrachten dazu einige Lehrbuch-Beispiele aus der Ökologie, siehe z.B.

 M. Begon, C. R. Townsend, and J. L. Harper.  
*Ecology: From Individuals to Ecosystems.*  
Blackell Publishing, 4 edition, 2008.

## Mittelwert und Standardabweichung. . .

- charakterisieren die Daten gut, falls deren Verteilung glockenförmig ist
- und müssen andernfalls mit Vorsicht interpretiert werden.

Wir betrachten dazu einige Lehrbuch-Beispiele aus der Ökologie, siehe z.B.

 M. Begon, C. R. Townsend, and J. L. Harper.  
*Ecology: From Individuals to Ecosystems.*  
Blackell Publishing, 4 edition, 2008.

Im Folgenden verwenden wir zum Teil simulierte Daten, wenn die Originaldaten nicht verfügbar waren.

(Nehmen Sie also nicht alle Datenpunkte wörtlich.)

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten**
  - Beispiel: Wählerische Bachstelzen**
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

# Bachstelzen fressen Dungfliegen

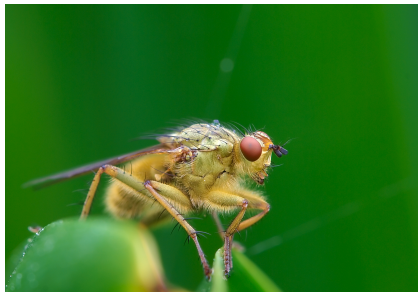
Räuber



Bachstelze (White Wagtail)  
*Motacilla alba alba*

image (c) by Artur Mikolajewski

Beute



Gelbe Dungfliege  
*Scatophaga stercoraria*

image (c) by Viatour Luc

# Vermutung

- Die Fliegen sind unterschiedlich groß
- Effizienz für die Bachstelze = Energiegewinn / Zeit zum Fangen und fressen
- Laborexperimente lassen vermuten, dass die Effizienz bei 7mm großen Fliegen maximal ist.

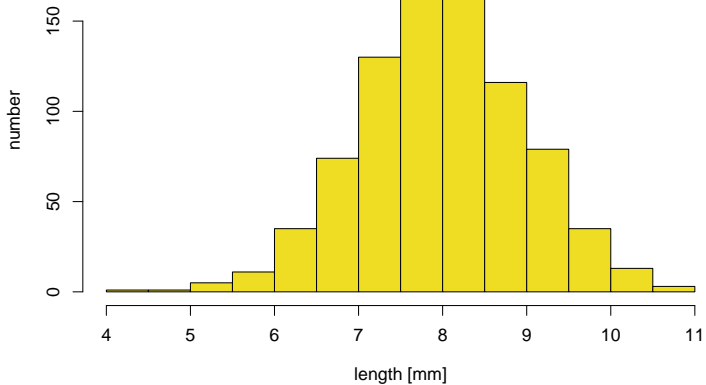


N.B. Davies.

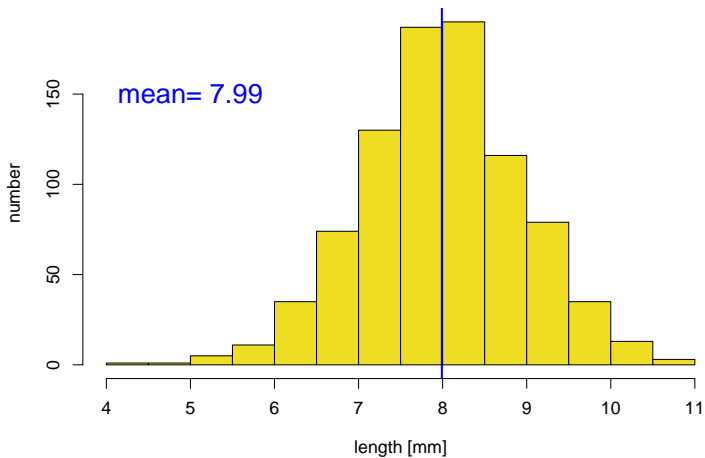
Prey selection and social behaviour in wagtails (Aves: Motacillidae).

*J. Anim. Ecol.*, 46:37–57, 1977.

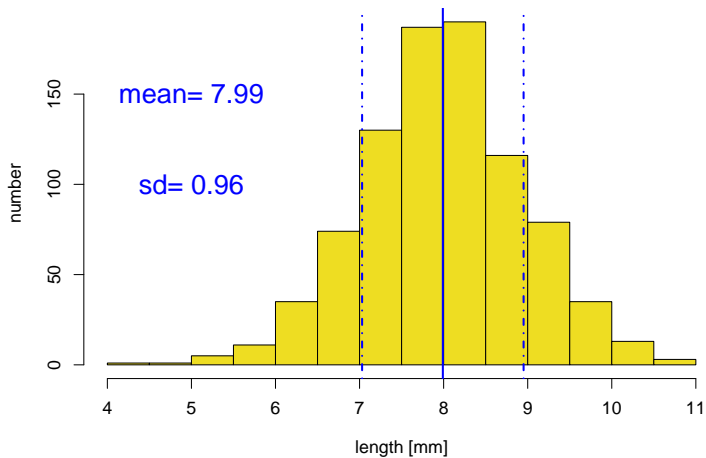
## available dung flies



## available dung flies

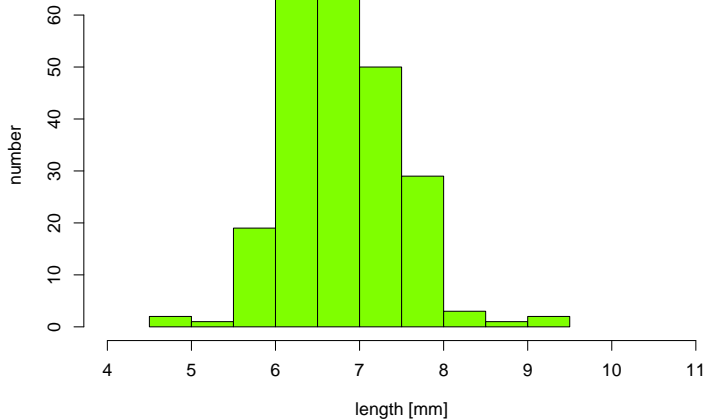


## available dung flies

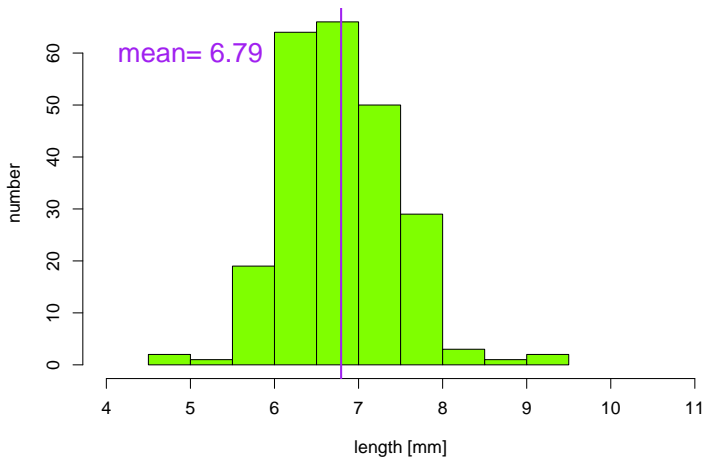




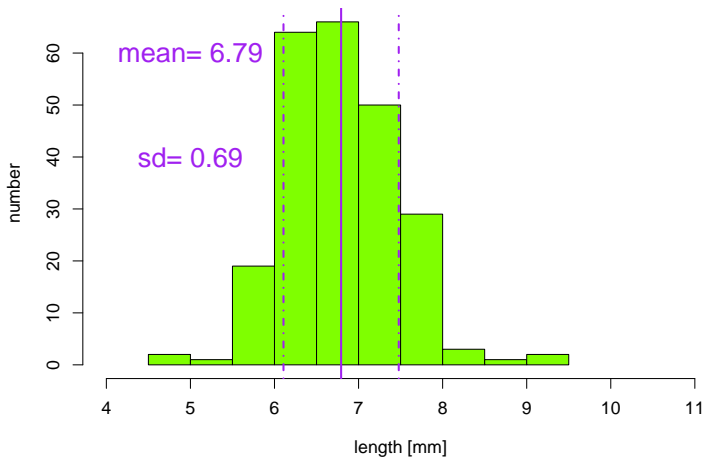
## captured dung flies



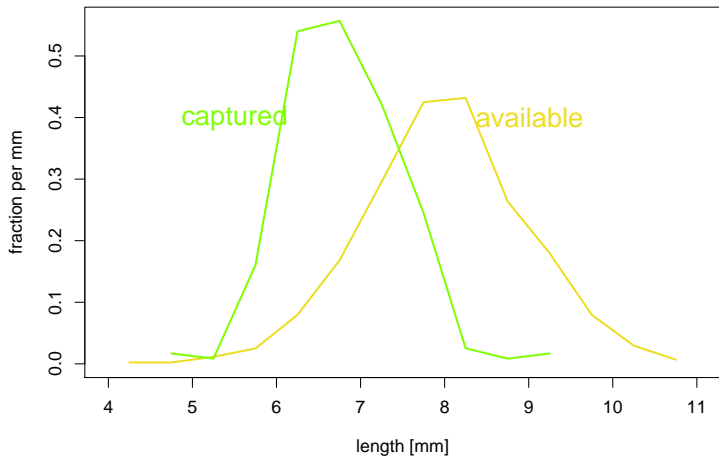
## captured dung flies



## captured dung flies



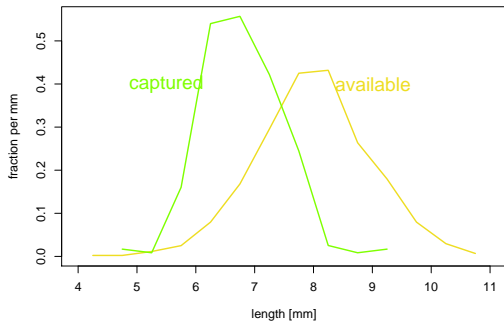
## dung flies: available, captured



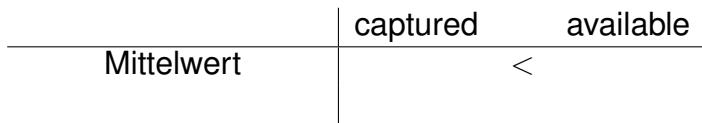
# Vergleich der Größenverteilungen

	captured	available
Mittelwert		

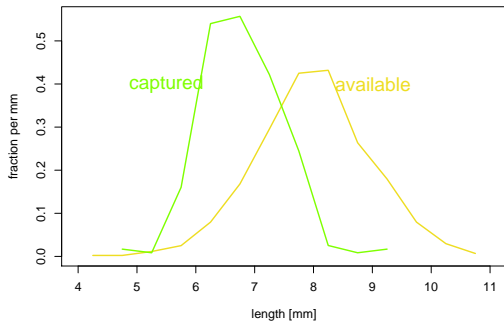
dung flies: available, captured



# Vergleich der Größenverteilungen



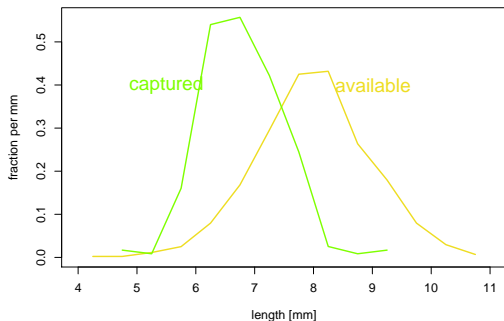
dung flies: available, captured



# Vergleich der Größenverteilungen

	captured		available
Mittelwert	6.29	<	7.99

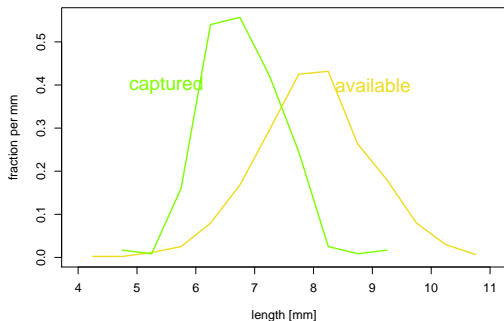
dung flies: available, captured



# Vergleich der Größenverteilungen

	captured	<	available
Mittelwert	6.29		7.99
Standardabweichung			

dung flies: available, captured

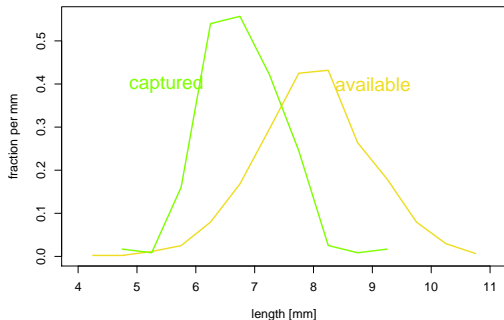




# Vergleich der Größenverteilungen

	captured		available
Mittelwert	6.29	<	7.99
Standardabweichung		<	

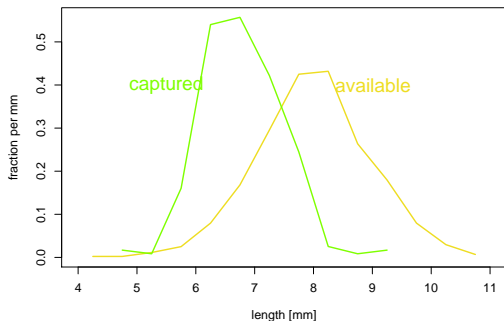
dung flies: available, captured



# Vergleich der Größenverteilungen

	captured		available
Mittelwert	6.29	<	7.99
Standardabweichung	0.69	<	0.96

dung flies: available, captured



# Interpretation

Die Bachstelzen bevorzugen Dungfliegen, die etwa 7mm groß sind.

# Interpretation

Die Bachstelzen bevorzugen Dungfliegen, die etwa 7mm groß sind.

Hier waren die Verteilungen glockenförmig und es genügten 4 Werte (die beiden Mittelwerte und die beiden Standardabweichungen), um die Daten adäquat zu beschreiben.

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten**
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman**
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R



*Nephila madagascariensis*

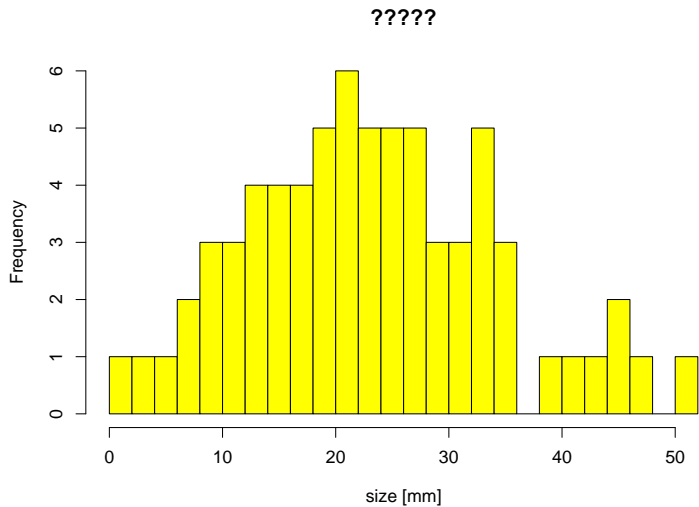
image (c) by Bernard Gagnon

## Simulierte Daten:

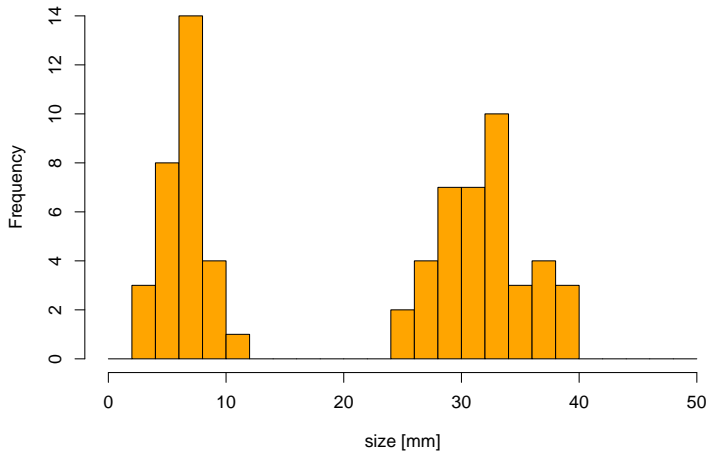
Eine Stichprobe von 70 Spinnen

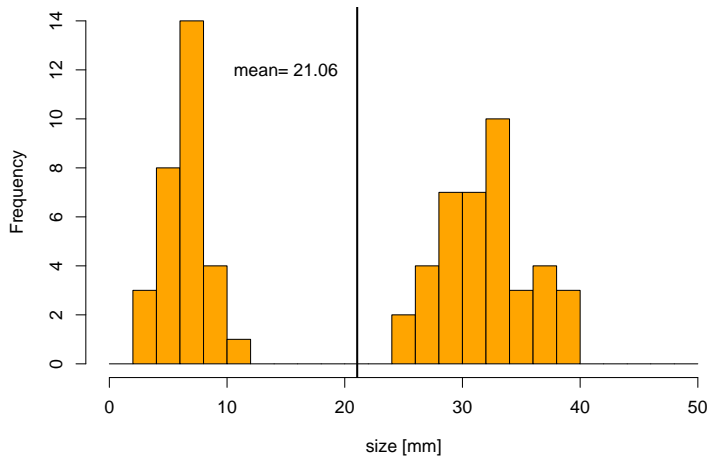
Mittlere Größe: 21,06 mm

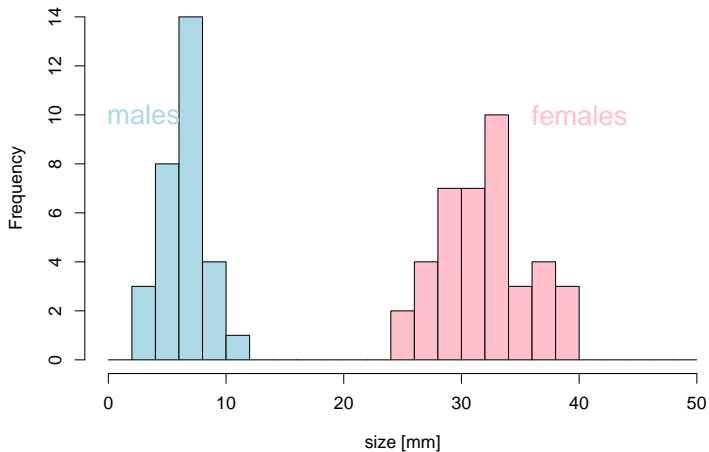
Standardabweichung der Größe: 12,94 mm

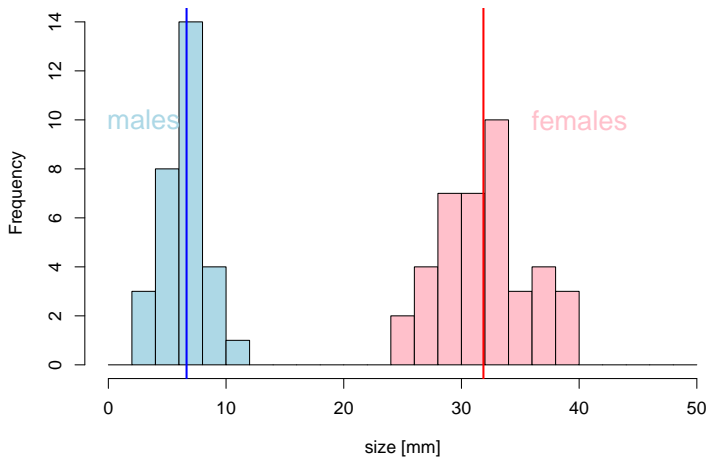




***Nephila madagascariensis* (n=70)**

***Nephila madagascariensis* (n=70)**

***Nephila madagascariensis* (n=70)**

***Nephila madagascariensis* (n=70)**



*Nephila madagascariensis*

image (c) by Arthur Chapman

# Fazit des Spinnenbeispiels

Wenn die Daten aus verschiedenen Gruppen zusammengesetzt sind, die sich bezüglich des Merkmals deutlich unterscheiden, kann es sinnvoll sein, Kenngrößen wie den Mittelwert für jede Gruppe einzeln zu berechnen.

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten**
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras**
- 5 Das Statistikpaket R
  - Deskriptive Statistik mit R

# Kupfertolerantes Rotes Straußgras



Rotes Straußgras  
*Agrostis tenuis*

image (c) Kristian Peters



Kupfer  
*Cuprum*

Hendrick met de Bles





A.D. Bradshaw.

Population Differentiation in *agrostis tenuis* Sibth. III.  
populations in varied environments.

*New Phytologist*, 59(1):92 – 103, 1960.



T. McNeilly and A.D Bradshaw.

Evolutionary Processes in Populations of Copper Tolerant  
*Agrostis tenuis* Sibth.

*Evolution*, 22:108–118, 1968.



A.D. Bradshaw.

Population Differentiation in *agrostis tenuis* Sibth. III.  
populations in varied environments.

*New Phytologist*, 59(1):92 – 103, 1960.



T. McNeilly and A.D Bradshaw.

Evolutionary Processes in Populations of Copper Tolerant  
*Agrostis tenuis* Sibth.

*Evolution*, 22:108–118, 1968.

Wir verwenden hier wieder simulierte Daten, da die  
Originaldaten nicht zur Verfügung stehen.

# Anpassung an Kupfer?

- Pflanzen, denen das Kupfer schadet, haben kürzere Wurzeln.

# Anpassung an Kupfer?

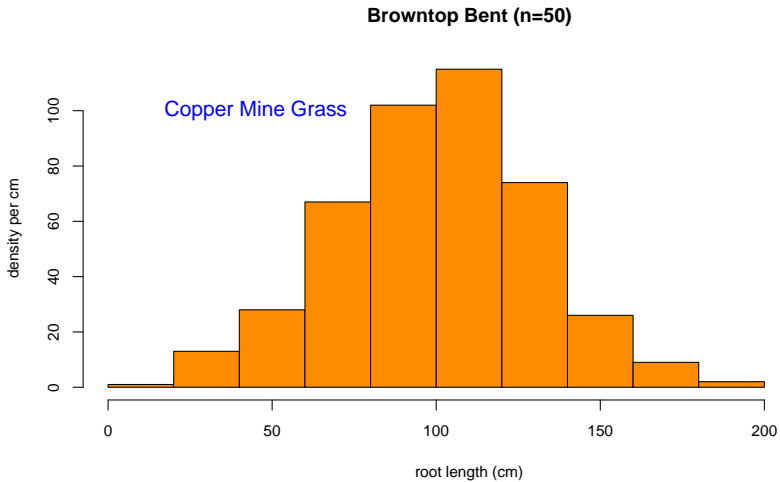
- Pflanzen, denen das Kupfer schadet, haben kürzere Wurzeln.
- Die Wurzellängen von Pflanzen aus der Umgebung von Kupferminen wird gemessen.

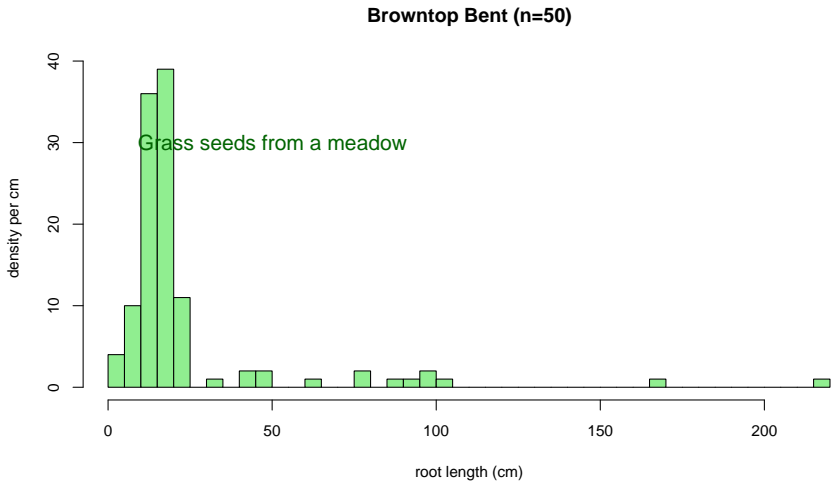
# Anpassung an Kupfer?

- Pflanzen, denen das Kupfer schadet, haben kürzere Wurzeln.
- Die Wurzellängen von Pflanzen aus der Umgebung von Kupferminen wird gemessen.
- Samen von unbelasteten Wiesen werden bei Kupferminen eingesät.

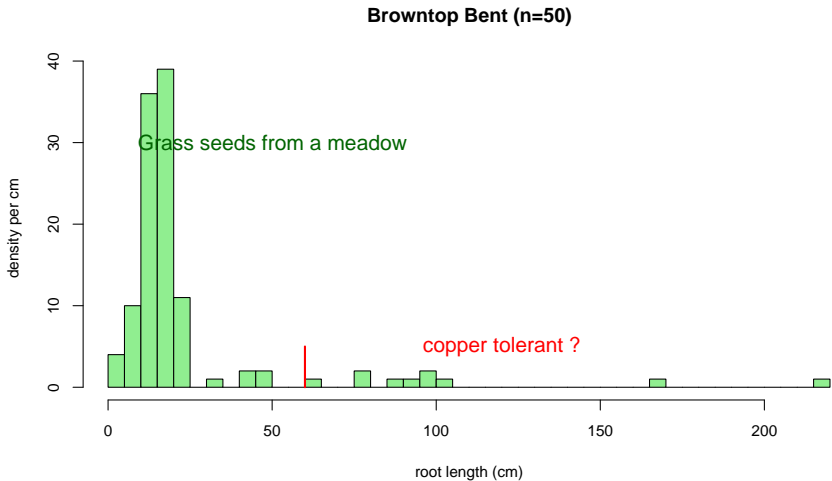
# Anpassung an Kupfer?

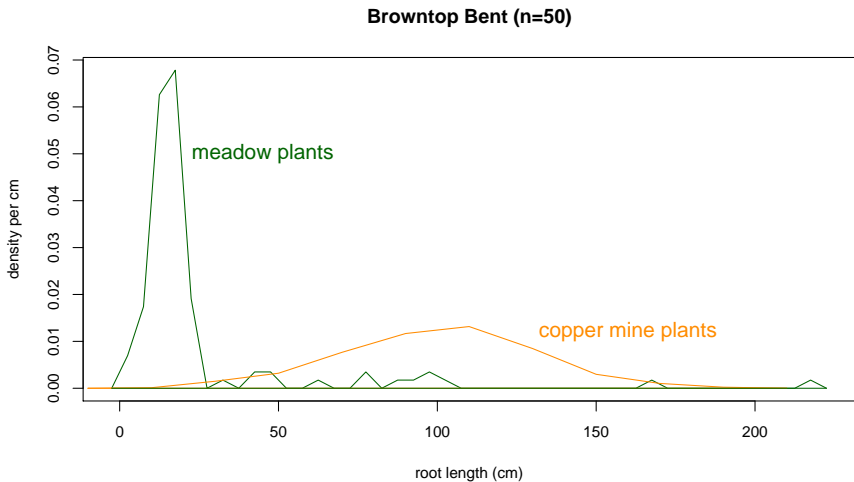
- Pflanzen, denen das Kupfer schadet, haben kürzere Wurzeln.
- Die Wurzellängen von Pflanzen aus der Umgebung von Kupferminen wird gemessen.
- Samen von unbelasteten Wiesen werden bei Kupferminen eingesät.
- Die Wurzellängen dieser “Wiesenpflanzen” werden gemessen.

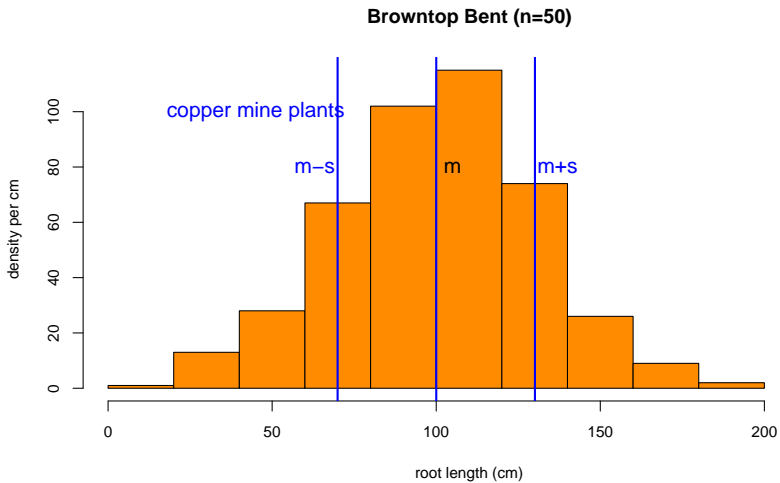


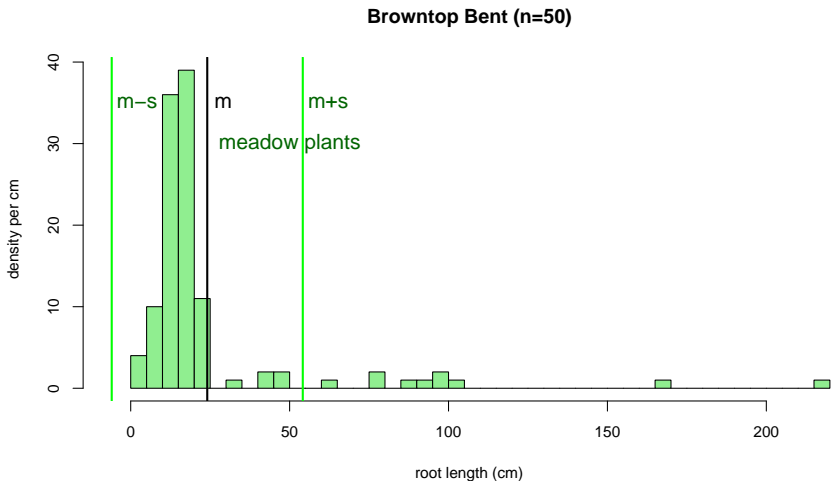












2/3 der Wurzellängen innerhalb  $[m-sd, m+sd]$ ???? **Nein!**

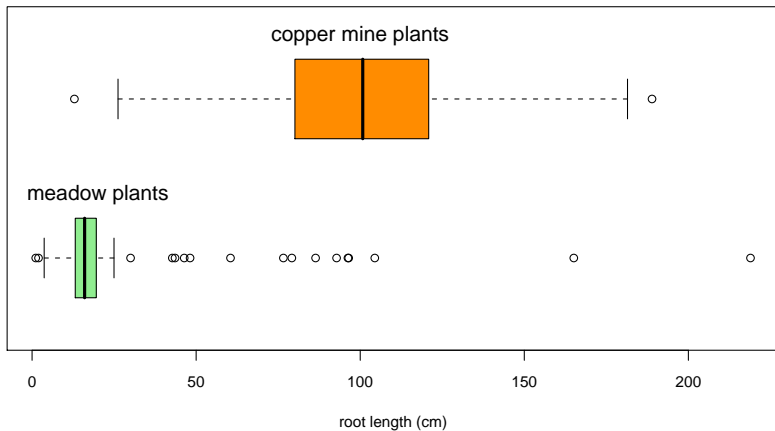
# Fazit des Straußgras-Beispiels

Manche Verteilungen können nur mit mehr als zwei Variablen angemessen beschrieben werden.

# Fazit des Straußgras-Beispiels

Manche Verteilungen können nur mit mehr als zwei Variablen angemessen beschrieben werden.

z.B. mit den fünf Werten der Boxplots:  
 $\min$ ,  $Q_1$ , median,  $Q_3$ ,  $\max$

**Browntop Bent n=50+50**

# Schlussfolgerung

Viele Datenverteilungen sind annähernd glockenförmig und können durch den **Mittelwert** und die **Standardabweichung** hinreichend beschrieben werden.



# Schlussfolgerung

Viele Datenverteilungen sind annähernd glockenförmig und können durch den **Mittelwert** und die **Standardabweichung** hinreichend beschrieben werden.

Es gibt aber auch Ausnahmen. Also: **Besser** ist es, die Daten auch graphisch zu untersuchen, und sich **nicht** allein auf numerische Kenngrößen zu verlassen.

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 **Das Statistikpaket R**
  - **Deskriptive Statistik mit R**

# Was ist R?

- Man kann R als (sehr mächtigen) „Statistik-Taschenrechner“ verwenden.

# Was ist R?

- Man kann R als (sehr mächtigen) „Statistik-Taschenrechner“ verwenden.

(R ist eine für die Statistik und für stochastische Simulation entwickelte Programmiersprache, zudem sind viele statistische Standardverfahren bereits in R implementiert oder als Zusatzpaket verfügbar.

R wird auch im Stochastik-Praktikums verwendet.)

# Was ist R?

- Man kann R als (sehr mächtigen) „Statistik-Taschenrechner“ verwenden.

(R ist eine für die Statistik und für stochastische Simulation entwickelte Programmiersprache, zudem sind viele statistische Standardverfahren bereits in R implementiert oder als Zusatzpaket verfügbar.

R wird auch im Stochastik-Praktikums verwendet.)

- R hat eine sehr aktive Benutzer- und Entwicklergemeinschaft (die nahezu alle Bereiche der Statistik und viele Anwendungsbereiche (z.B. Populationsgenetik, Finanzmathematik) überdeckt).

# Was ist R?

- Man kann R als (sehr mächtigen) „Statistik-Taschenrechner“ verwenden.

(R ist eine für die Statistik und für stochastische Simulation entwickelte Programmiersprache, zudem sind viele statistische Standardverfahren bereits in R implementiert oder als Zusatzpaket verfügbar.

R wird auch im Stochastik-Praktikums verwendet.)

- R hat eine sehr aktive Benutzer- und Entwicklergemeinschaft (die nahezu alle Bereiche der Statistik und viele Anwendungsbereiche (z.B. Populationsgenetik, Finanzmathematik) überdeckt).
- R ist frei verfügbar unter der GNU general public license, für (nahezu) alle Rechnerarchitekturen erhältlich:  
<http://www.r-project.org/>
- R ist auf ZDV-Rechnern installiert.

# R installieren, starten, anhalten

Installation: Windows, Mac OS: Binaries von

<http://www.r-project.org/> (siehe Link Download, Packages, CRAN dort)

Linux: Für die meisten Distributionen gibt es fertige Pakete

Fragen oder Probleme:

Tutorium (J. Blauth), Do 10-12, Raum MI 2

R starten: Windows, Mac OS: Icon (ggf. aus Menu) anklicken,

Linux/Unix: `>` R auf einer Konsole

(oder mit ESS aus Emacs heraus, oder aus `rstudio`, ...).

R beenden: `q()` (fragt, ob Daten gespeichert werden sollen)

laufende Rechnungen unterbrechen: `CTRL-C`

# Inhalt

- 1 Ansatz der Statistik
- 2 Graphische Darstellungen
  - Histogramme und Dichtepolygone
  - Stripcharts
  - Boxplots
- 3 Statistische Kenngrößen
  - Median und andere Quartile, (empirische) Quantile
  - Mittelwert und Standardabweichung
- 4 Beispiele zum Sinn und Unsinn von Mittelwerten
  - Beispiel: Wählerische Bachstelzen
  - Beispiel: Spiderman & Spiderwoman
  - Beispiel: Kupfertoleranz beim Roten Straußgras
- 5 **Das Statistikpaket R**
  - **Deskriptive Statistik mit R**



## Datensatz `x` in R eingeben

```
x <- c( 53,52,41,41,42,58,40,43,42,38,43,49,34,51,45,  
       39,41,45,45,39,37,36,42,44,47,43,46,43,43,45,  
       42,52,49,44,50,40,47,46,50,50,41,51,41,47,42,  
       52,36,46,42,56,39,40,36,42,36,36,47,45,47,49 )
```

(aus Datei einlesen: `x <- scan('Dateiname')`)

Mittelwert (`mean`), Standardabweichung (`sd`), Median, und

Quantile

```
mean(x)
```

```
sd(x)
```

```
median(x)
```

```
quantile(x, 0.25, type=1)
```

```
quantile(x, 0.75, type=1)
```

```
summary(x)
```

Boxplot, Histogramm

```
boxplot(x)
```

```
hist(x)
```

## Nur zur Information: Literatur zu R

Wir werden im Zusammenhang der Vorlesung nur (fakultativ) einige wenige R-Befehle ansprechen, so dass Sie für etwaige freiwillige „R-Schnupperei“ i.A. keine Literatur benötigen werden.


- „Standardreferenz“: W.N. Venables et al, *An Introduction to R*,  
<http://cran.r-project.org/manuals.html>
- Günther Sawitzki, *Einführung in R*,  
<http://sintro.r-forge.r-project.org/>
- William N. Venables, Brian D. Ripley, *Modern applied statistics with S*  
(„Standardlehrbuch“, UB Lehrbuchsammlung)
- Lothar Sachs and Jürgen Hedderich, *Angewandte Statistik – Methodensammlung mit R* (E-Book, UB)
- Christine Duller, *Einführung in die nichtparametrische Statistik mit SAS und R : ein anwendungsorientiertes Lehr- und Arbeitsbuch* (E-Book, UB)
- Helge Toutenburg, Christian Heumann, *Deskriptive Statistik : Eine Einführung in Methoden und Anwendungen mit R und SPSS* (E-Book, UB)
- Uwe Ligges, *Programmieren mit R* (E-Book, UB)

The R Project for Statistical Computing - Mozilla Firefox

The R Project for Statistical Co...  
www.r-project.org

Meistbesucht Getting Started Latest Headlines

## The R Project for Statistical Computing




About R  
[What is R?](#)  
[Contributors](#)  
[Screenshots](#)  
[What's new?](#)

Download, Packages  
[CRAN](#)

R Project  
[Foundation](#)  
[Members & Donors](#)  
[Mailing Lists](#)  
[Bug Tracking](#)  
[Developer Page](#)  
[Conferences](#)  
[Search](#)

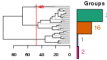
Documentation  
[Manuals](#)  
[FAQs](#)  
[The R Journal](#)  
[Wiki](#)  
[Books](#)  
[Certification](#)  
[Other](#)

Misc  
[Bioconductor](#)  
[Related Projects](#)  
[User Groups](#)  
[Links](#)




PCA 5 vars  
`plot(pcaobj[,c("dim1", "dim2")])`

Clustering: 4 groups



Group: 28, 16, 2



Getting Started:

- R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To [download R](#), please choose your preferred [CRAN mirror](#).
- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

News:

- R version 3.0.2** (Frisbee Sailing) has been released on 2013-09-25.
- useR! 2013**, took place at the University of Castilla-La Mancha, Albacete, Spain, July 10-12 2013.
- The R Journal Vol.5/1** is available.
- R version 2.15.3** (Security Blanket) has been released on 2013-03-01.

This server is hosted by the [Institute for Statistics and Mathematics](#) of [WU \(Wirtschaftsuniversität Wien\)](#).

# Nochmal: Idee der Statistik

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

# Nochmal: Idee der Statistik

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

Die Daten werden als Realisierungen von Zufallsvariablen  
aufgefasst, die in einem stochastischen Modell spezifiziert  
werden.

# Nochmal: Idee der Statistik

Variabilität (Erscheinung der Natur) durch Zufall  
(mathematische Abstraktion) modellieren

Die Daten werden als Realisierungen von Zufallsvariablen  
aufgefasst, die in einem stochastischen Modell spezifiziert  
werden.

Man versucht dann, anhand der Daten Rückschlüsse auf  
Parameter des Modells zu ziehen, und so systematische Effekte  
von Zufälligem zu trennen.