

# Biostatistik, WS 2010/2011

## Deskriptive Statistik

Matthias Birkner

<http://www.mathematik.uni-mainz.de/~birkner/Biostatistik1011/>

19.11.2010



1/112

Wozu Statistik?

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

Andrejs Dunkels

3/112

## Was ist Statistik?

Die Natur ist voller Variabilität.

Wie geht man mit variablen Daten um?

Es gibt eine mathematische Theorie des

Zufalls:

die **Stochastik**.

## IDEE DER STATISTIK

**Variabilität**

(Erscheinung der Natur)

durch

**Zufall**

(mathematische Abstraktion)

**modellieren.**

Statistik

=

Datenanalyse  
mit Hilfe  
stochastischer Modelle

## Beispiel

Daten aus einer Diplomarbeit aus 2001 am  
Forschungsinstitut Senckenberg, Frankfurt  
am Main

Crustaceensektion

Leitung: *Dr. Michael Türkay*



*Charybdis acutidens* TÜRKAY 1985

## Der Springkrebs *Galathea intermedia*



9/112

## Helgoländer Tiefe Rinne, Fang vom 6.9.1988

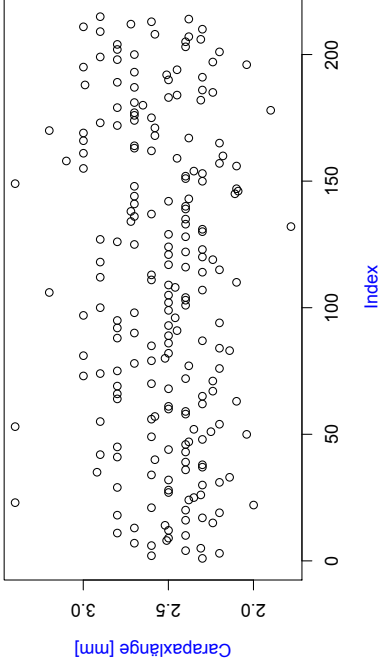
### Carapaxlänge (mm):

#### Nichteiertragende Weibchen ( $n = 215$ )

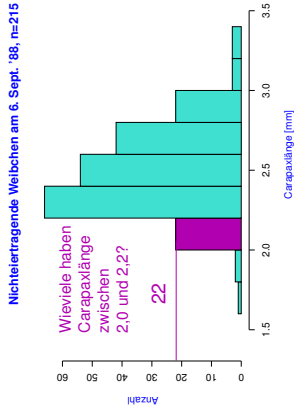
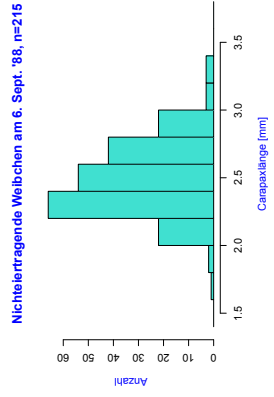
2,9	3,0	2,9	2,5	2,7	2,9	2,9	3,0
3,0	2,9	3,4	2,8	2,9	2,8	2,8	2,4
2,8	2,5	2,7	3,0	2,9	3,2	3,1	3,0
2,7	2,5	3,0	2,8	2,8	2,8	2,7	3,0
2,6	3,0	2,9	2,8	2,9	2,9	2,3	2,7
2,6	2,7	2,5	.	.	.	.	.

10/112

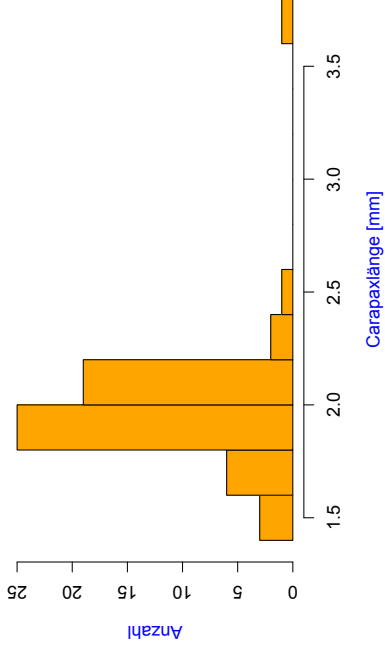
**Nichteiertragende Weibchen am 6. Sept. '88, n=215**



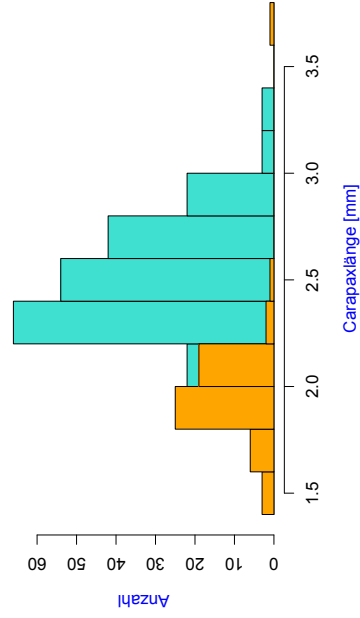
Eine Möglichkeit der graphischen  
Darstellung:  
das Histogramm



Analoge Daten zwei Monate später  
(3.11.88):

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

16/112

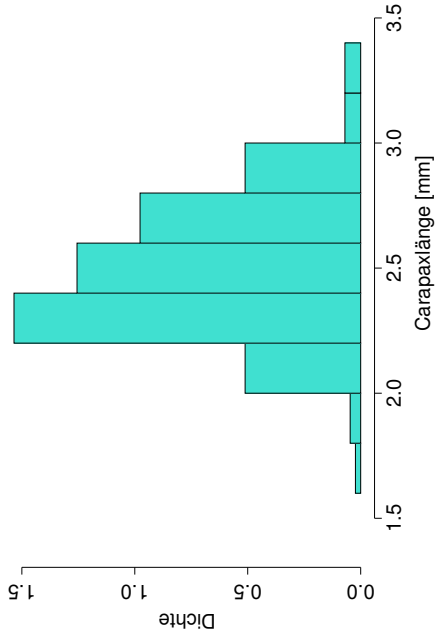
**Vergleich der beiden Verteilungen****Nichteiertragende Weibchen**

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

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

17/112

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

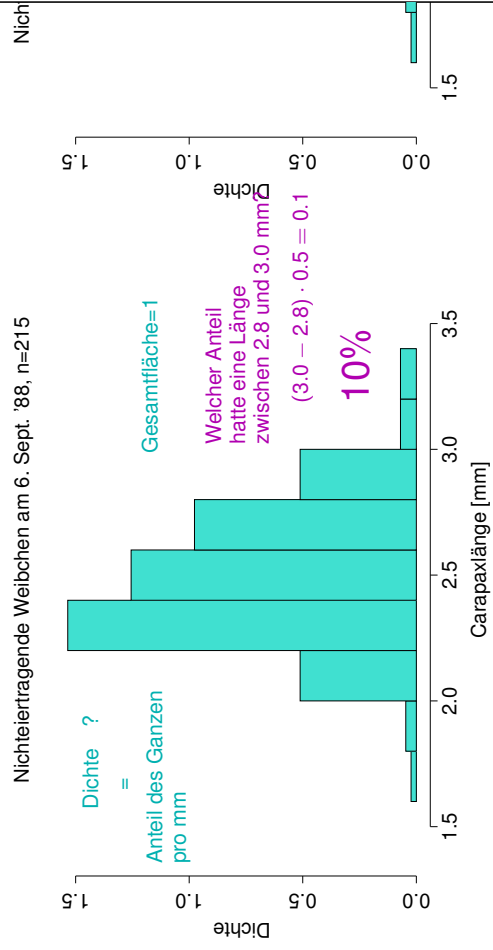


18/112

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

19/112



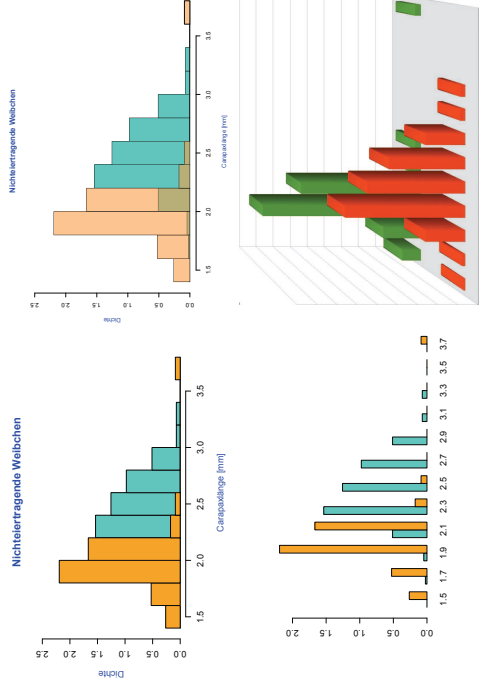


20/112

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

21/112

Versuche, die Histogramme zusammen zu zeigen:



22/112

## Ratschlag

Wenn Sie Schauwerbegestalter(in) sind:

Beeindrucken Sie Jung und Alt mit total abgefahrener 3D-Plots!

Wenn Sie Wissenschaftler(in) werden wollen:

Bevorzugen Sie einfache und klare 2D-Darstellungen.

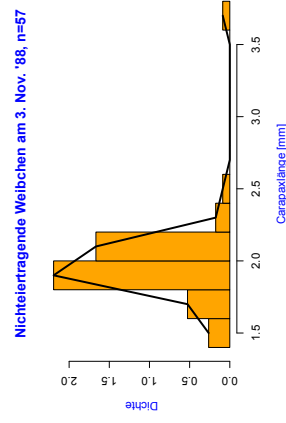
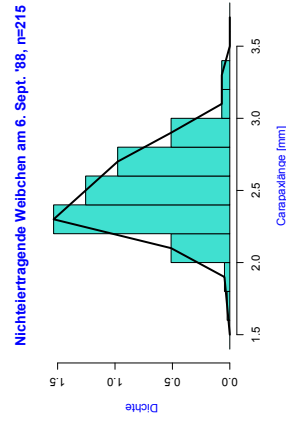
23/112

## Problem

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

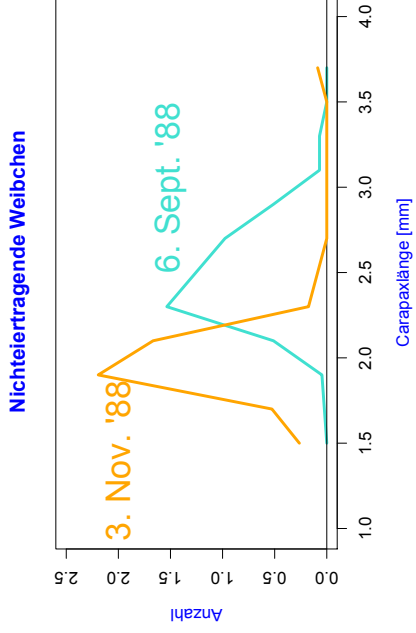
24/112

## Einfache und klare Lösung: Dichtepolygone



25/112

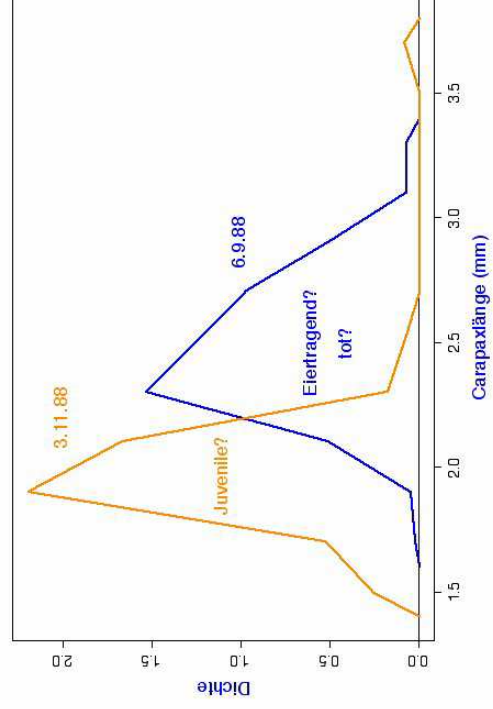
## Zwei und mehr Dichtepolygone in einem Plot



Biologische Interpretation der Verschiebung?

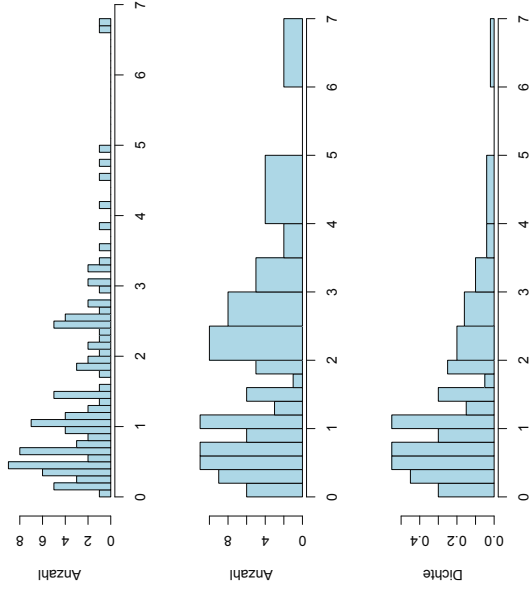
26/112

## Nichteiertragende Weibchen 6.9.88 und 3.11.88

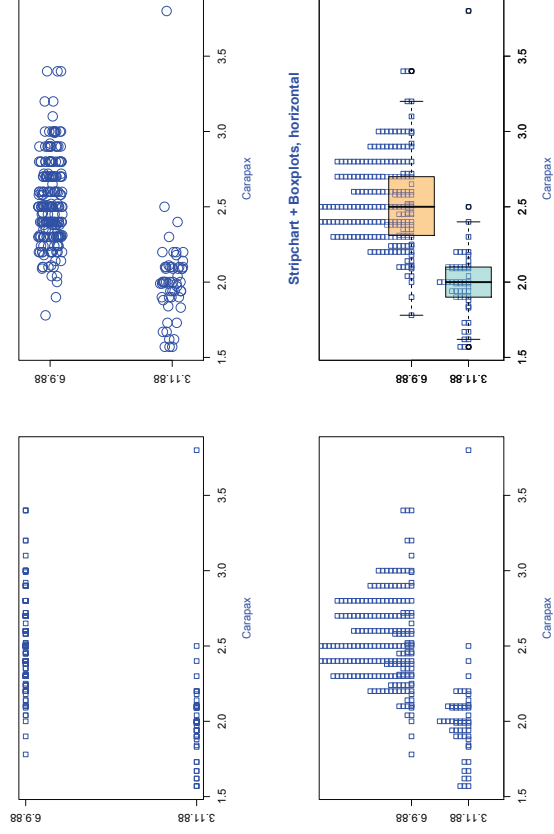


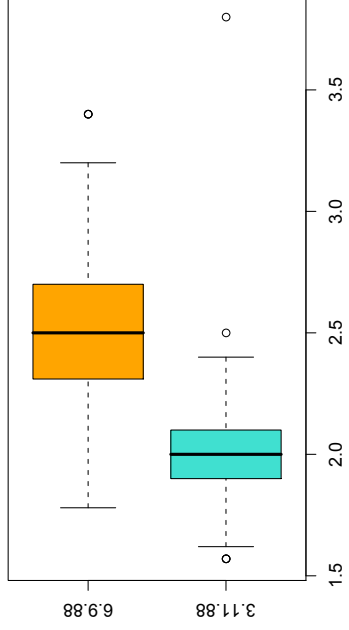
27/112

# Anzahl vs. Dichte



Also: Bei Histogrammen mit ungleichmäßiger Unterteilung immer Dichten verwenden!



**Boxplots, horizontal**

31/112

Histogramme und Dichtepolygone  
geben  
ein ausführliches Bild  
eines Datensatzes.

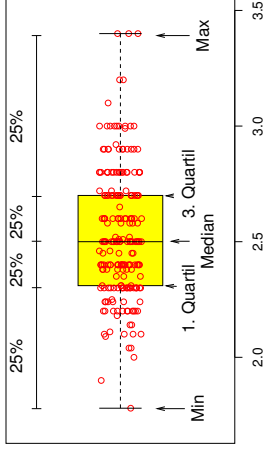
**Manchmal zu ausführlich.**

32/112

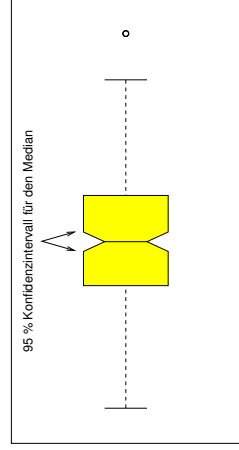


# Der Boxplot

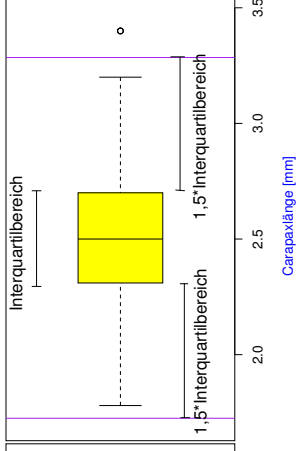
Boxplot, einfache Ausführung



Boxplot, Profiausstattung



Boxplot, Standardausführung



36/112

Beispiel:  
Die Ringeltaube  
*Palumbus palumbus*

38/112

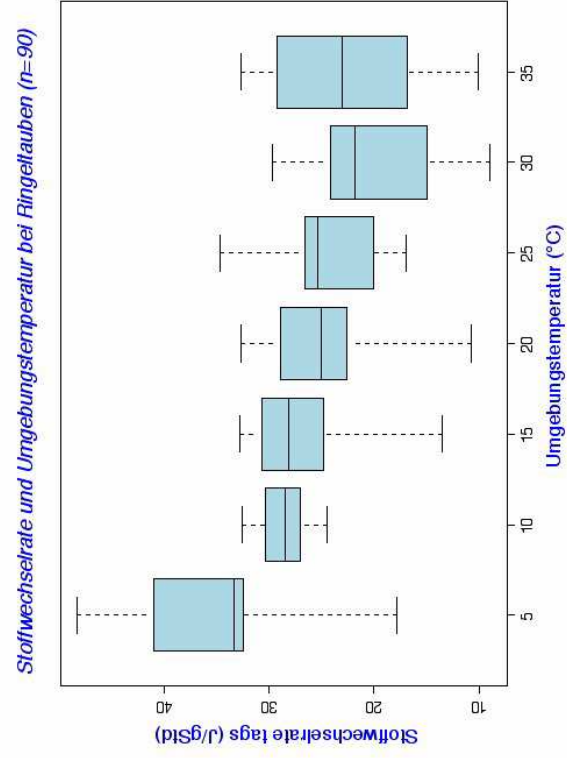




Wie hängt die Stoffwechselrate bei der Ringeltaube von der Umgebungstemperatur ab?

# Daten aus dem AK Stoffwechselphysiologie Prof. Prinzinger Universität Frankfurt

41/112

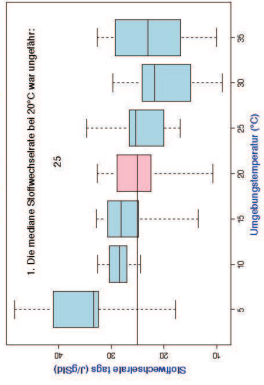


42/112

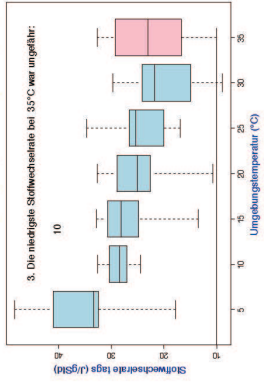
Klar:  
Stoffwechselrate  
höher  
bei  
tiefen Temperaturen

Vermutung:  
Bei hohen Temperaturen  
nimmt die Stoffwechselrate  
wieder zu  
(Hitzestress).

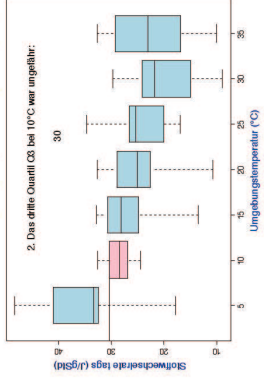
Stoffwechsellage und Umgebungstemperatur bei Ringeltauben (n=20)



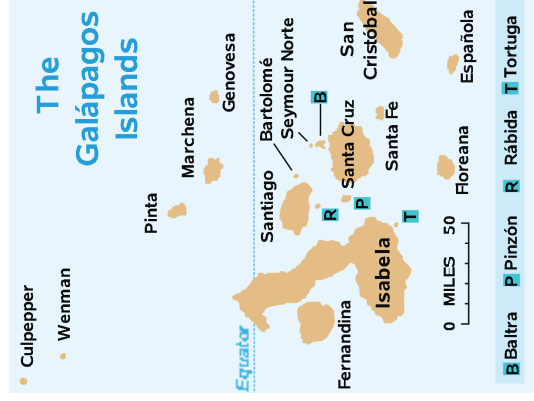
Stoffwechsellage und Umgebungstemperatur bei Ringeltauben (n=20)



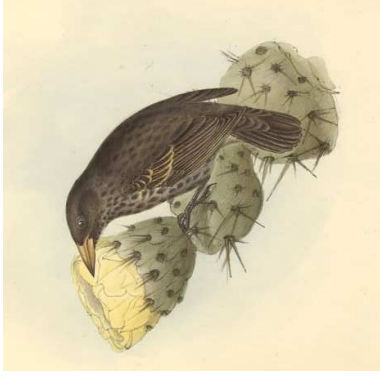
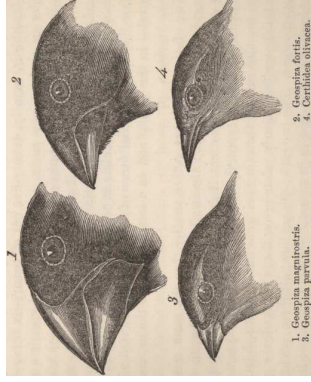
Stoffwechsellage und Umgebungstemperatur bei Ringeltauben (n=20)



# Charles Robert Darwin (1809-1882)



## Darwin-Finken




[http:](http://darwin-online.org.uk/graphics/Zoology_Illustrations.html)

[//darwin-online.org.uk/graphics/Zoology\\_Illustrations.html](http://darwin-online.org.uk/graphics/Zoology_Illustrations.html)

48/112

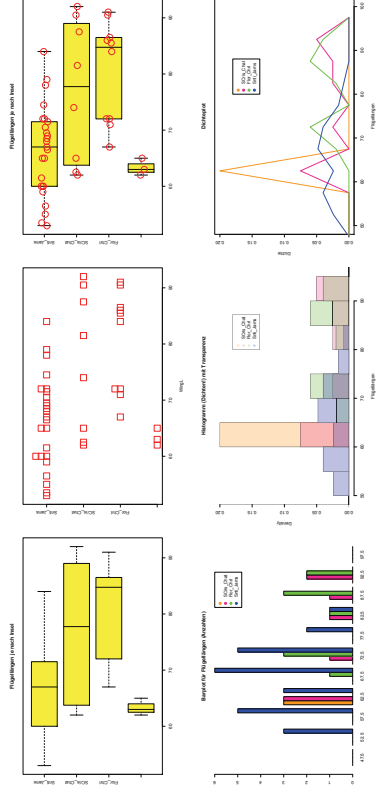
## Darwins Finken-Sammlung

 Sulloway, F.J. (1982) The Beagle collections of Darwin's Finches (Geospizinae). *Bulletin of the British Museum (Natural History)*, *Zoology series* **43**: 49-94.

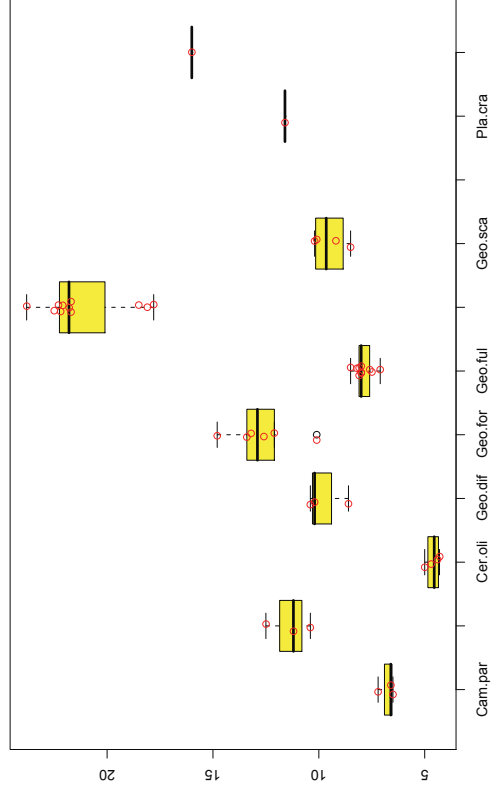
▶ <http://datadryad.org/repo/handle/10255/dryad.154>

49/112

# Flügelängen der Darwin-Finken



## Schnabelgröße je nach Art



## 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 verwenden
- 5 Vorsicht mit Tricks wie 3D oder halbttransparenten Farben
- 6 Jeder Datensatz ist anders; keine Patentrezepte

52/112

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

54/112

Wesentlich:

1. Wie groß?

Lageparameter

2. Wie variabel?

Streuungsparameter

Eine Möglichkeit  
kennen wir schon  
aus dem Boxplot:



Lageparameter

Der Median

Streuungsparameter

Der Quartilabstand ( $Q_3 - Q_1$ )

Der Median:  
die Hälfte der Beobachtungen sind kleiner,  
die Hälfte sind größer.

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

## Die Quartile

**Das erste Quartil,  $Q_1$ :**  
ein Viertel der Beobachtungen  
sind kleiner,  
drei Viertel sind größer.

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

60/112

## Die Quartile

**Das dritte Quartil,  $Q_3$ :**  
drei Viertel der Beobachtungen  
sind kleiner,  
ein Viertel sind größer.

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

61/112

Am häufigsten werden benutzt:

Lageparameter

Der Mittelwert  $\bar{x}$

Streuungsparameter

Die Standardabweichung  $s$

Der Mittelwert  
(engl. *mean*)

**NOTATION:**

Wenn die Beobachtungen

$x_1, x_2, x_3, \dots, x_n$

heißen,

schreibt man oft

$\bar{x}$

für den Mittelwert.

65/112

**DEFINITION:**

Mittelwert =  
$$\frac{\text{Summe der Messwerte}}{\text{Anzahl der Messwerte}}$$

Der Mittelwert von  $x_1, x_2, \dots, x_n$  als Formel:

$$\begin{aligned}\bar{x} &= (x_1 + x_2 + \dots + x_n) / n \\ &= \frac{1}{n} \sum_{i=1}^n x_i\end{aligned}$$

66/112

Beispiel:

$$x_1 = 3, x_2 = 0, x_3 = 2, x_4 = 3, x_5 = 1$$

$$\bar{x} = \text{Summe} / \text{Anzahl}$$

$$\bar{x} = (3 + 0 + 2 + 3 + 1) / 5$$

$$\bar{x} = 9 / 5$$

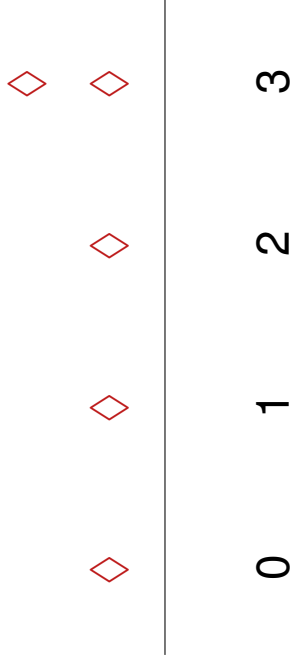
$$\bar{x} = 1,8$$

Geometrische Bedeutung  
des Mittelwerts:  
Der Schwerpunkt

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

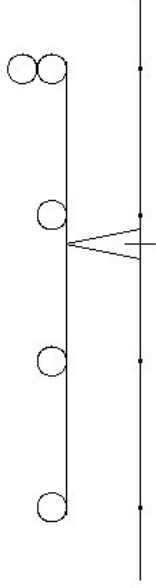
vor:

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



69/112

$m = 1,8$  ?



richtig

70/112

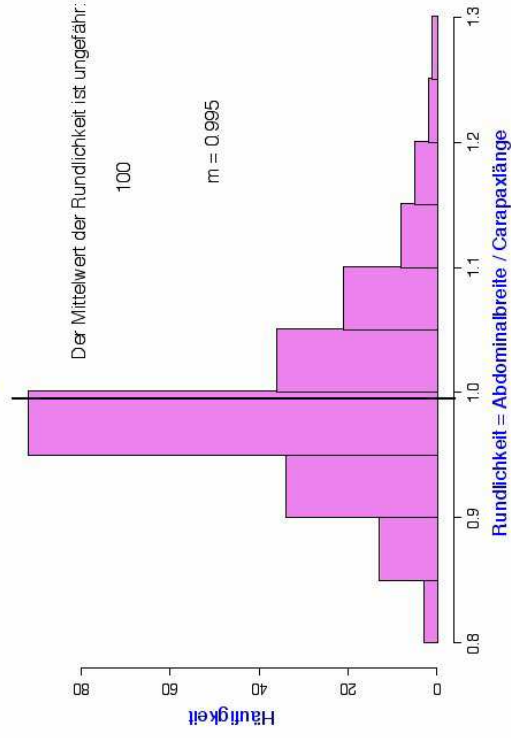
Beispiel: *Galathea intermedia*

„Rundlichkeit“  
:=  
Abdominalbreite / Carapaxlänge

Vermutung:  
Rundlichkeit nimmt  
bei Geschlechtsreife zu

71/112

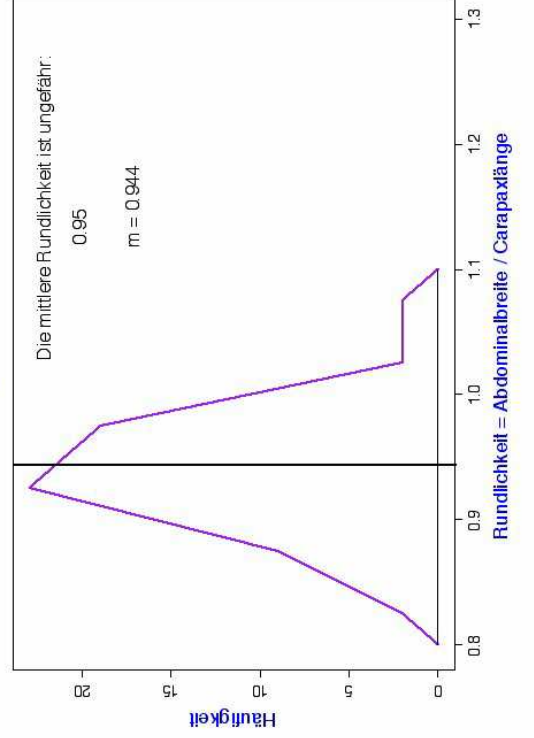
*Nichttragende Weibchen 6.9.88*



72/112

# Beispiel: 3.11.88

## Nichtertragende Weibchen 3.11.88



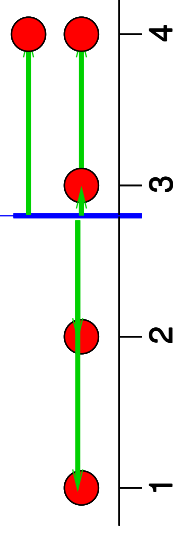


Die Standardabweichung  
Wie weit weicht  
eine typische Beobachtung  
vom  
Mittelwert  
ab ?

75/112

typische Mittelwert=2,8

Abweichung = ~~2,8~~ - 2,88 = -0,08



76/112

Die **Standardabweichung**  $\sigma$  (“sigma”) [auch *SD* von engl. *standard deviation*] ist ein

etwas komisches  
gewichtetes Mittel  
der Abweichungsbeträge  
und zwar

$$\sigma = \sqrt{\text{Summe}(\text{Abweichungen}^2)/n}$$

77/112

Die **Standardabweichung** von  $x_1, x_2, \dots, x_n$  als Formel:

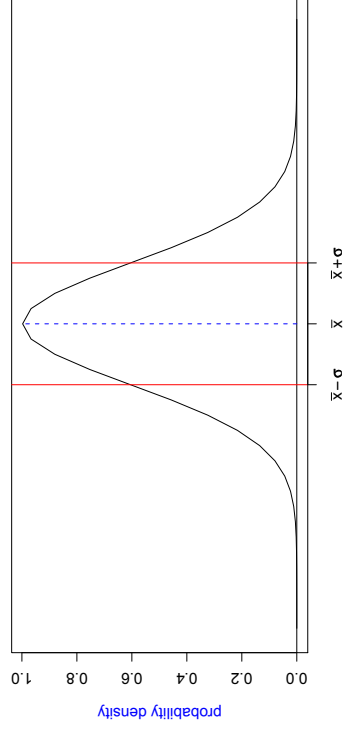
$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

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

78/112

## Faustregel für die Standardabweichung

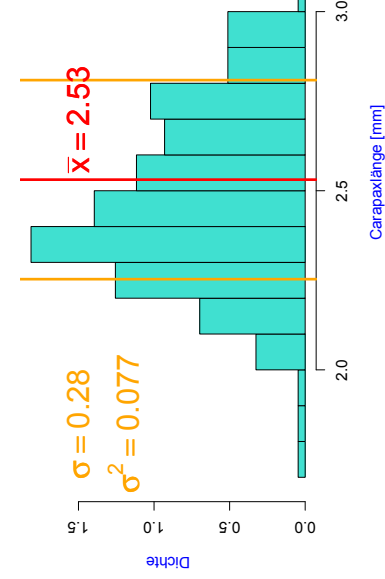
Bei ungefähr glockenförmigen (also eingipfligen und symmetrischen) Verteilungen liegen ca. 2/3 der Verteilung zwischen  $\bar{x} - \sigma$  und  $\bar{x} + \sigma$ .



79/112

## Standardabweichung der Carapaxlängen nichteiertragender Weibchen vom 6.9.88

### Nichteiertragende Weibchen



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

80/112

## Varianz der Carapaxlängen nichtleitertragender Weibchen vom 6.9.88

Alle Carapaxlängen im Meer:  $\mathcal{X} = (X_1, X_2, \dots, X_N)$ .  
Carapaxlängen in unserer Stichprobe:  $\mathcal{S} = (S_1, S_2, \dots, S_{n=215})$   
Stichprobenvarianz:

$$\sigma_S^2 = \frac{1}{n} \sum_{i=1}^{215} (s_i - \bar{s})^2 \approx 0,0768$$

Können wir 0,0768 als Schätzwert für die Varianz  $\sigma_X^2$  in der ganzen Population verwenden?

Ja, können wir machen. Allerdings ist  $\sigma_S^2$  im Durchschnitt um den Faktor  $\frac{n-1}{n}$  (= 214/215  $\approx$  0,995) kleiner als  $\sigma_X^2$

81/112

## Varianzbegriffe

Varianz in der Population:  $\sigma_X^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$

Stichprobenvarianz:  $\sigma_S^2 = \frac{1}{n} \sum_{i=1}^n (S_i - \bar{S})^2$   
korrigierte Stichprobenvarianz:

$$\begin{aligned} s^2 &= \frac{n}{n-1} \sigma_S^2 \\ &= \frac{n}{n-1} \cdot \frac{1}{n} \cdot \sum_{i=1}^n (s_i - \bar{s})^2 \\ &= \frac{1}{n-1} \cdot \sum_{i=1}^n (s_i - \bar{s})^2 \end{aligned}$$

Mit "Standardabweichung von  $S$ " ist meistens das korrigierte  $s$  gemeint.

82/112

$$\bar{x} = 10/5 = 2$$

Summe

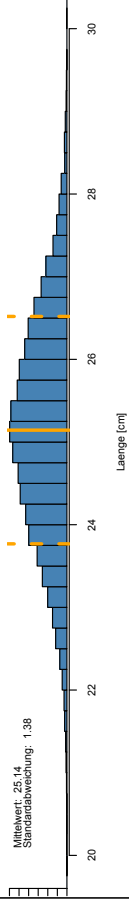
x	1	3	0	5	1	10
$x - \bar{x}$	-1	1	-2	3	-1	0
$(x - \bar{x})^2$	1	1	4	9	1	16

$$s^2 = \text{Summe}((x - \bar{x})^2) / (n - 1)$$

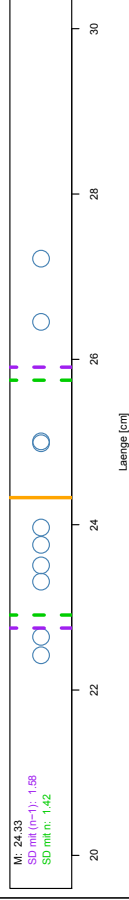
$$= 16 / (5 - 1) = 4$$

$$s = 2$$

Eine simulierte Fischpopulation (N=10000 adulte)

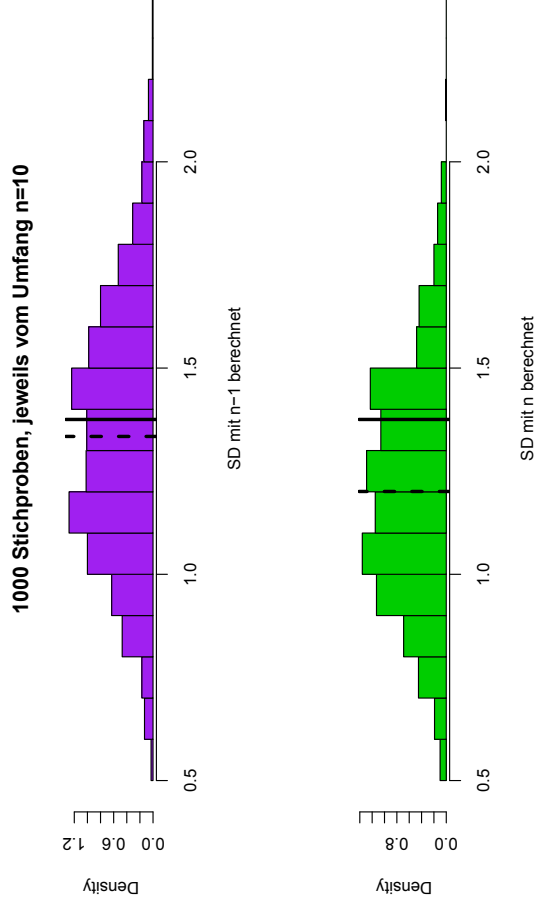


Eine Stichprobe aus der Population (n=10)



Noch eine Stichprobe aus der Population (n=10)





85/112

## $\sigma$ mit $n$ oder $n - 1$ berechnen?

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

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

Wenn es sich bei  $x_1, \dots, x_n$  um eine Stichprobe handelt (wie meistens in der Statistik), sollten Sie die Formel

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

verwenden.

86/112

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. Glauben Sie uns also nicht alle Datenpunkte.**

88/112

## Bachstelzen fressen Dungfliegen

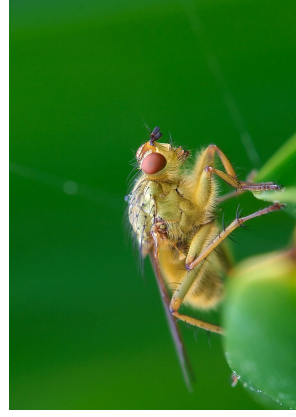
Räuber



Bachstelze (White Wagtail)  
*Motacilla alba*

image (c) by Artur Mikolajewski

Beute



Gelbe Dungfliege  
*Scatophaga stercoraria*

image (c) by Viatour Luc

90/112

## 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.

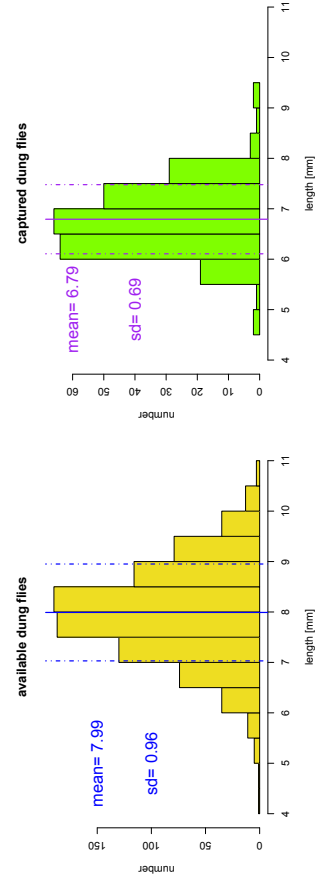


**N.B. Davies.**

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

*J. Anim. Ecol.*, 46:37–57, 1977.

91/112

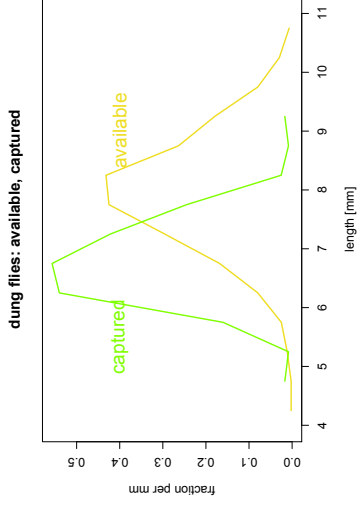


92/112



## Vergleich der Größenverteilungen

	captured	available
Mittelwert	6.29	< 7.99
Standardabweichung	0.69	< 0.96



93/112

## 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.

94/112



*Nephila madagascariensis*

image (c) by Bernard Gagnon

96/112

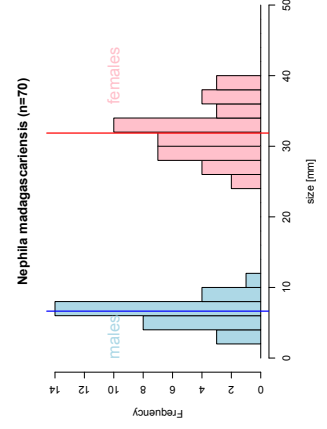
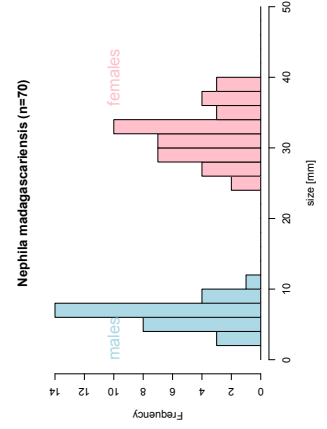
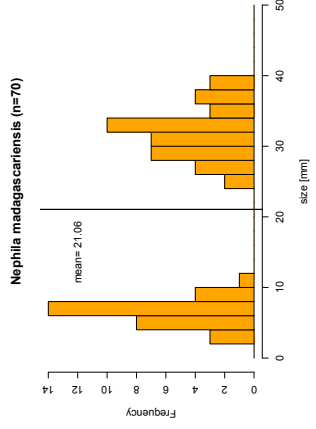
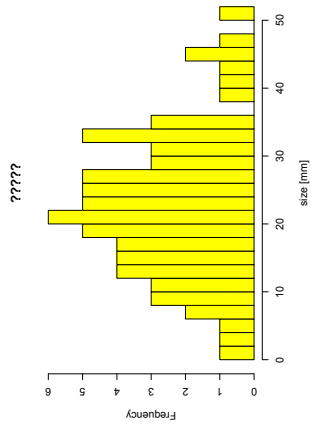
**Simulated Data:**

Eine Stichprobe von 70 Spinnen

Mittlere Größe: 21,06 mm

Standardabweichung der Größe: 12,94 mm

97/112



*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.

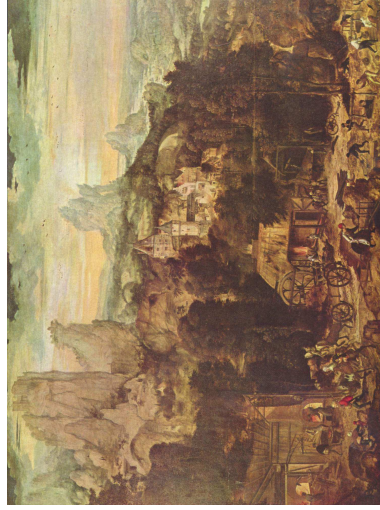
100/112

## Kupfertolerantes Rotes Straußgras



Rotes Straußgras  
*Agrostis tenuis*

image (c) Kristian Peters



Kupfer  
*Cuprum*

Hendrick met de Bles

102/112

 A.D. Bradshaw.

Population Differentiation in *agrostis tenuis* Sibth. III. populations in varied environments. *New Phytologist*, 59(1):92 – 103, 1960.

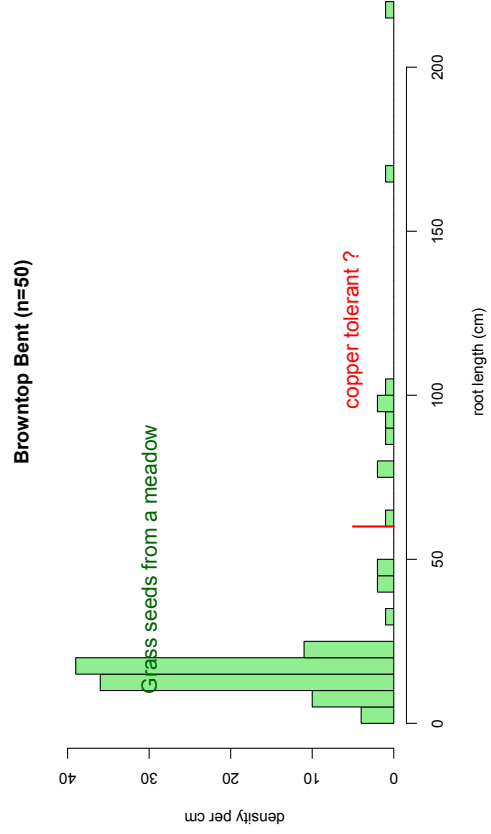
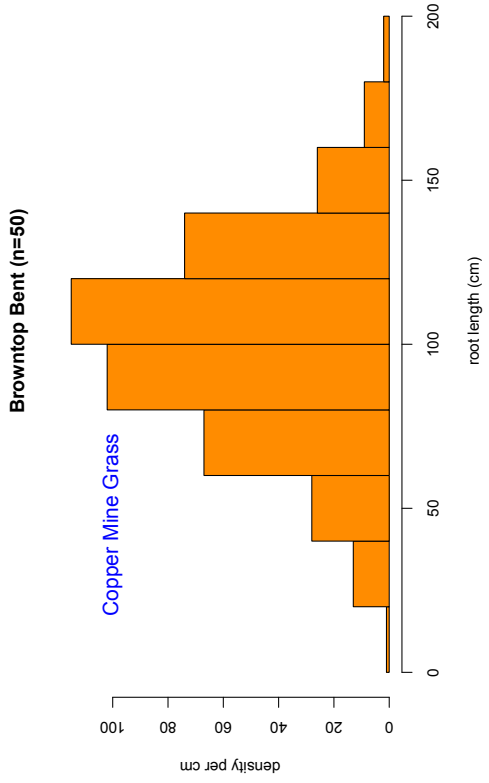
 T. McNeilly and A.D Bradshaw.

Evolutionary Processes in Populations of Copper Tolerant *Agrostis tenuis* Sibth. *Evolution*, 22:108–118, 1968.

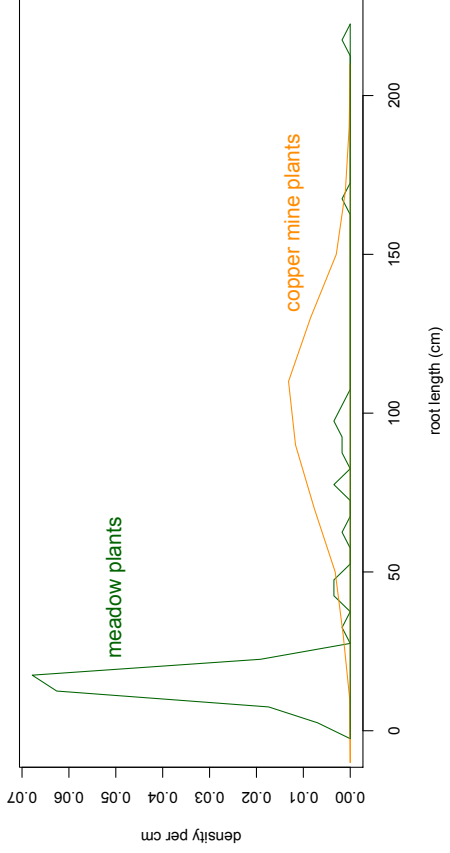
Wir verwenden hier wieder simulierte Daten, da die Originaldaten nicht zur Verfügung stehen.

## Anpassung an Kupfer?

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

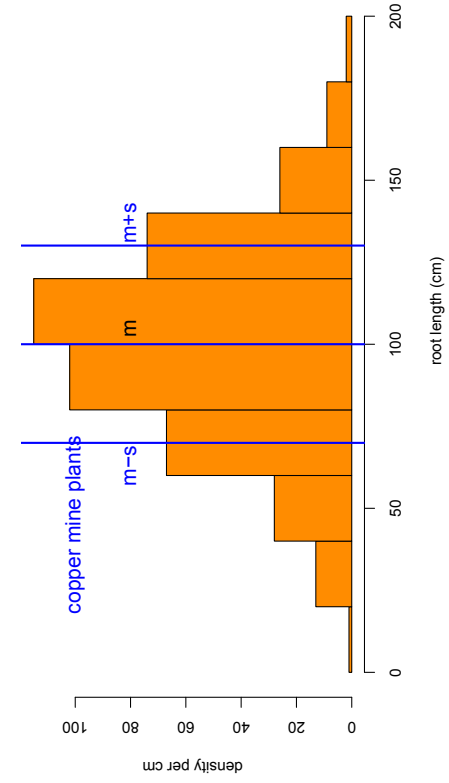


### Browntop Bent (n=50)

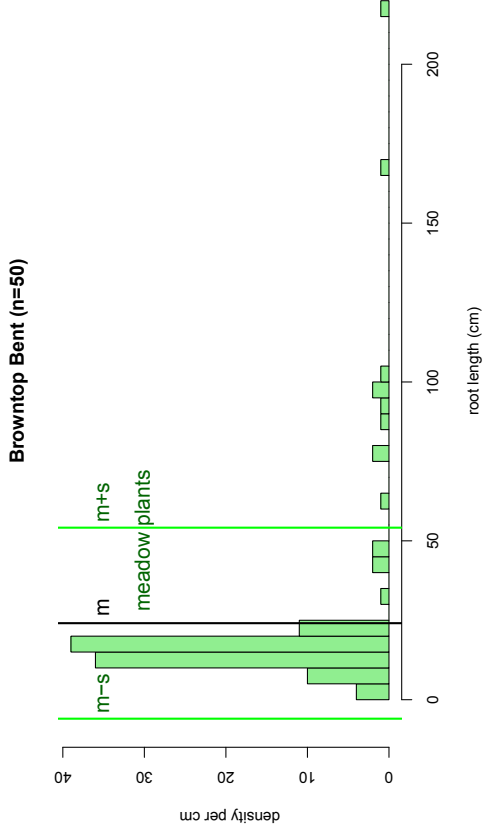


107/112

### Browntop Bent (n=50)



108/112



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

109/112

## Fazit des Straußgras-Beispiels

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

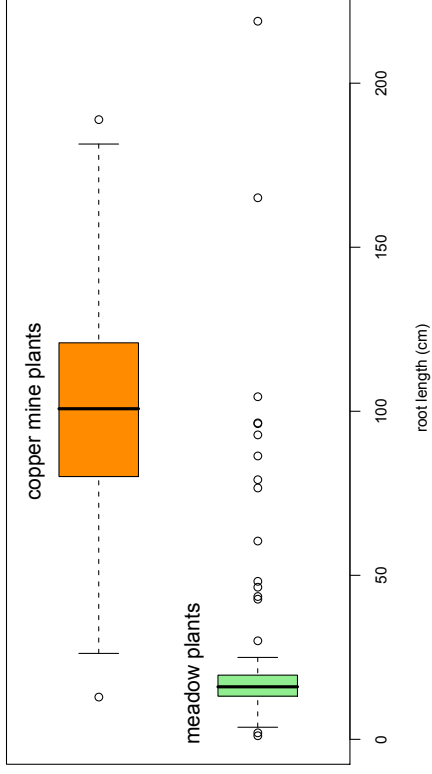
z.B. mit den fünf Werten der Boxplots:

min,  $Q_1$ , median,  $Q_3$ , max

110/112



## Browntop Bent n=50+50



111/112

## Schlussfolgerung

In der Biologie sind viele Datenverteilungen annähernd glockenförmig und können durch den Mittelwert und die Standardabweichung hinreichend beschrieben werden.

Es gibt aber auch Ausnahmen. Also:  
**Immer** die Daten erst mal graphisch untersuchen!  
Verlassen sie sich **niemals** allein auf numerische Kenngrößen!

112/112