

# Biostatistik, WS 2013/2014

## Lineare Regression

Matthias Birkner

<http://www.mathematik.uni-mainz.de/~birkner/Biostatistik1314/>

24.1.2014



JOHANNES GUTENBERG  
UNIVERSITÄT MAINZ



photo (c) by Jörg Hempel

*Gypus fulvus*  
Gänsegeier

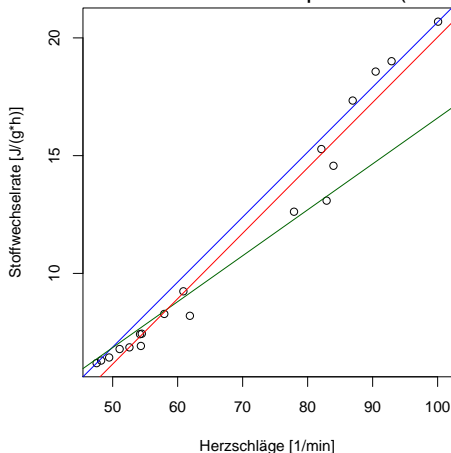


Prinzinger, R., E. Karl, R. Bögel, Ch. Walzer (1999): Energy metabolism, body temperature, and cardiac work in the Griffon vulture *Gyps vulvus* - telemetric investigations in the laboratory and in the field.  
*Zoology* **102**, Suppl. II: 15

- Daten aus der Arbeitsgruppe Stoffwechselphysiologie (Prof. Prinzinger) der Frankfurter Goethe-Universität.
- Telemetrisches System zur Messung der Herzfrequenz bei Vögeln auch während des Fluges.
- Wichtig für ökologische Fragen: die Stoffwechselrate
- Messung der Stoffwechselrate aufwändig und nur im Labor möglich.
- Können wir von der Herzfrequenz auf die Stoffwechselrate schließen?

Die Daten:

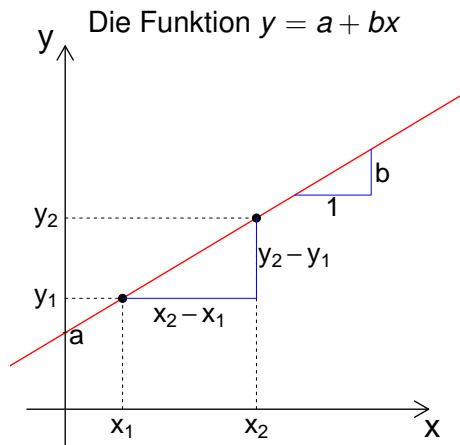
19 Beobachtungen eines Geiers  
im Labor bei konstanter Temperatur (16° C)



Die Beobachtungen legen einen linearen Zusammenhang zwischen Herzfrequenz und Stoffwechselrate nahe.

Welche Gerade passt „am Besten“?

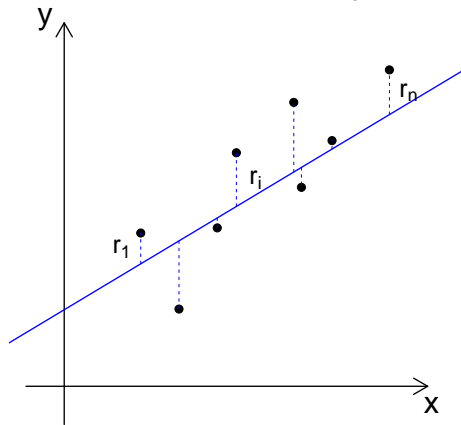
## Erinnerung: Lineare Funktionen



Hier ist  $a$  der Achsenabschnitt und  $b$  die Steigung.  
 Für jedes Paar von Punkten  $(x_1, y_1) \neq (x_2, y_2)$  auf der Geraden  
 gilt  $\frac{y_2 - y_1}{x_2 - x_1} = b$ .

Unsere Situation:  $n$  beobachtete Paare  $(x_i, y_i)$ ,  $i = 1, \dots, n$

Die Daten und eine Gerade  $y = a + bx$

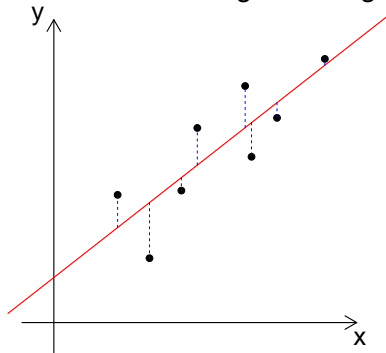


Für keine gegebene Steigung  $b$  und Achsenabschnitt  $a$  liegen die Beobachtungen genau auf der Geraden  $y = a + bx$ .

Es bleiben stets sogenannte *Residuen*  $r_i := y_i - a - bx_i (\neq 0)$ .

# kleinste-Quadrate-Schätzer $\hat{a}$ , $\hat{b}$ und Regressionsgerade

Die Daten und die Regressionsgerade



Ansatz: Finde  $\hat{a}$  und  $\hat{b}$  derart, dass  $\sum_{i=1}^n (y_i - \hat{a} - \hat{b} x_i)^2 \stackrel{!}{=} \text{minimal}$   
(wobei über alle Wahlen von  $a, b \in \mathbb{R}$  minimiert wird).

Die Gerade  $y = \hat{a} + \hat{b} x$  heißt die *Regressionsgerade*.

Man kann  $\hat{a}$ ,  $\hat{b}$  folgendermaßen berechnen:

Seien  $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\bar{y} := \frac{1}{n} \sum_{i=1}^n y_i$  die (empirischen) Mittelwerte der  $x$ - bzw. der  $y$ -Werte,

$\sigma_x^2 := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ ,  $\sigma_y^2 := \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$ , die zugehörigen („unkorrigierten Stichproben“-)Varianzen und

(„unkorrigierten Stichproben“-)Varianzen und

$\text{cov}(x, y) := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$  die (empirische) Kovarianz der  $x$ - und der  $y$ -Werte.

### Satz (Koeffizienten der Regressionsgerade)

$$\hat{b} = \frac{\text{cov}(x, y)}{\sigma_x^2}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}.$$



## Bemerkungen zu den Regressionskoeffizienten

$$\hat{b} = \frac{\text{cov}(x, y)}{\sigma_x^2}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

- 1 Beobachtung: Die Regressionsgerade  $y = \hat{a} + \hat{b}x$  geht durch den Schwerpunkt der Daten  $(\bar{x}, \bar{y})$ .  
(Dies kann eine Merkhilfe für die Formeln bilden.)

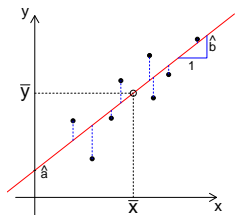
- 2 Man kann  $\hat{b}$  auch anders ausdrücken:

$$\frac{\text{cov}(x, y)}{\sigma_x^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n y_i(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

(Für das 1. Gleichheitszeichen erweitere mit  $n$ , für das 2. beachte  $\bar{y} \sum_{i=1}^n (x_i - \bar{x}) = 0$ . Die varianten Formeln können ggf. zum Rechnen angenehmer sein.)

## Bemerkungen zu den Regressionskoeffizienten

$$\hat{b} = \frac{\text{cov}(x, y)}{\sigma_x^2}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$



- 1 Beobachtung: Die Regressionsgerade  $y = \hat{a} + \hat{b}x$  geht durch den Schwerpunkt der Daten  $(\bar{x}, \bar{y})$ .  
(Dies kann eine Merkhilfe für die Formeln bilden.)

- 2 Man kann  $\hat{b}$  auch anders ausdrücken:

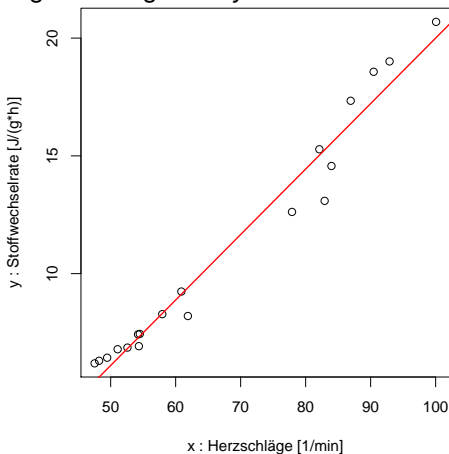
$$\frac{\text{cov}(x, y)}{\sigma_x^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n y_i(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

(Für das 1. Gleichheitszeichen erweitere mit  $n$ , für das 2. beachte  $\bar{y} \sum_{i=1}^n (x_i - \bar{x}) = 0$ . Die varianten Formeln können ggf. zum Rechnen angenehmer sein.)

Für das Geier-Beispiel ist mit  $x \hat{=}$  Herzfrequenz,  $y \hat{=}$  Stoffwechselrate:  
 $\bar{x} = 67.9$ ,  $\bar{y} = 11.1$ ,  $\text{cov}(x, y) = 83.5$ ,  $\sigma_x^2 = 300.7$ , also

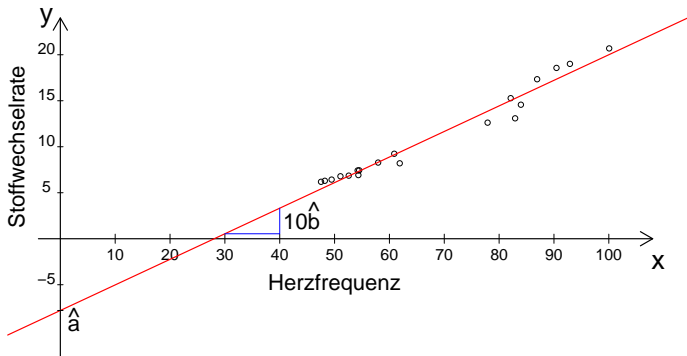
$$\hat{b} = \frac{83.5}{300.7} = 0.278, \quad \hat{a} = 11.1 - 67.9 \cdot 0.278 = -7.8$$

$n = 19$  Datenpunkte und  
 Regressionsgerade  $y = -7.8 + 0.278x$



## Interpretation der Regressionskoeffizienten

$n = 19$  Datenpunkte und Regressionsgerade  $y = -7.8 + 0.278x$



$\hat{b} = 0.278$  : Erhöhung der Herzfrequenz um 10 erhöht die Stoffwechselrate im Mittel um  $10\hat{b} = 2.78$ .

$\hat{a} = -7.8$  : Dies wäre die Stoffwechselrate eines hypothetischen Geiers mit Herzfrequenz 0.

(Offensichtlich kein sinnvoller Wert: Die Regressionsgerade ist nur in dem Bereich plausibel, in dem tatsächlich Beobachtungen vorliegen; Extrapolation auf eigene Gefahr!)

## Regressionskoeffizienten: Woher kommen die Formeln?

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n r_i^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - a - bx_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y} - b(x_i - \bar{x}))^2 + (\bar{y} - a - b\bar{x})^2 \end{aligned}$$

denn

$$\begin{aligned} (y_i - \bar{y} - b(x_i - \bar{x}))^2 &= (y_i - a - bx_i - (\bar{y} - a - b\bar{x}))^2 \\ &= (y_i - a - bx_i)^2 - 2(y_i - a - bx_i)(\bar{y} - a - b\bar{x}) + (\bar{y} - a - b\bar{x})^2 \end{aligned}$$

und wenn man über  $i$  summiert und durch  $n$  teilt, ergeben die beiden letzten Terme zusammen gerade  $-(\bar{y} - a - b\bar{x})^2$ .

## Regressionskoeffizienten: Woher kommen die Formeln?

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n r_i^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - a - bx_i)^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y} - b(x_i - \bar{x}))^2 + (\bar{y} - a - b\bar{x})^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 - 2b \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\
&\quad + b^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{y} - a - b\bar{x})^2 \\
&= \sigma_y^2 - 2b \operatorname{cov}(x, y) + b^2 \sigma_x^2 + (\bar{y} - a - b\bar{x})^2 \\
&= \sigma_y^2 - \frac{(\operatorname{cov}(x, y))^2}{\sigma_x^2} + \sigma_x^2 \left( b - \frac{\operatorname{cov}(x, y)}{\sigma_x^2} \right)^2 + (\bar{y} - a - b\bar{x})^2
\end{aligned}$$

## Regressionskoeffizienten: Woher kommen die Formeln?

Demnach:

$$\frac{1}{n} \sum_{i=1}^n r_i^2$$

$$= \underbrace{\sigma_y^2 - \frac{(\text{Cov}(x, y))^2}{\sigma_x^2}}_{\substack{\text{hängt nicht von } a, b \\ \text{ab (und ist } \geq 0)}} + \underbrace{\sigma_x^2 \left( b - \frac{\text{Cov}(x, y)}{\sigma_x^2} \right)^2}_{\geq 0} + \underbrace{(\bar{y} - a - b\bar{x})^2}_{\geq 0}$$

Die Summe der Residuenquadrate wird also minimiert durch die Wahlen

$$\hat{b} = \frac{\text{Cov}(x, y)}{\sigma_x^2}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}.$$

(Denn dann sind die beiden letzten Terme oben = 0.)

## Lineare Regression: Modellvorstellung

Wir haben die Regressionsgerade

$$y = \hat{a} + \hat{b} \cdot x$$

durch die Minimierung der Summe der quadrierten Residuen definiert:

$$(\hat{a}, \hat{b}) = \arg \min_{(a,b)} \sum_i (y_i - (a + b \cdot x_i))^2$$

Dahinter steckt die Modellvorstellung, dass Werte  $a, b$  existieren, so dass für alle Datenpaare  $(x_i, y_i)$  gilt

$$y_i = a + b \cdot x_i + \varepsilon_i,$$

wobei alle  $\varepsilon_i$  unabhängig und normalverteilt sind mit Mittelwert 0 und derselben Varianz  $\sigma^2$ .



gegebene Daten:

<b>Y</b>	<b>X</b>
$y_1$	$x_1$
$y_2$	$x_2$
$y_3$	$x_3$
$\vdots$	$\vdots$
$y_n$	$x_n$

Modell: es gibt Zahlen  
 $a, b, \sigma^2$ , so dass

$$\begin{aligned}
 y_1 &= a + b \cdot x_1 + \varepsilon_1 \\
 y_2 &= a + b \cdot x_2 + \varepsilon_2 \\
 y_3 &= a + b \cdot x_3 + \varepsilon_3 \\
 &\vdots \\
 y_n &= a + b \cdot x_n + \varepsilon_n
 \end{aligned}$$

Dabei sind  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  unabhängig  $\sim \mathcal{N}(0, \sigma^2)$ .

$\Rightarrow y_1, y_2, \dots, y_n$  sind unabhängig  $y_i \sim \mathcal{N}(a + b \cdot x_i, \sigma^2)$ .

$a, b, \sigma^2$  sind unbekannt, aber **nicht zufällig**.

## Lineare Regression: Theorie

Modell:  $y_i = a + b \cdot x_i + \varepsilon_i$  mit  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  
 fest (aber z.T. unbekannt):  $a, b, x_i, \sigma^2$  zufällig:  $\varepsilon_i, y_i$

Die kleinste-Quadrate-Schätzer

$$\hat{b} = \frac{\sum_i y_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}, \quad \hat{a} = \bar{y} - b\bar{x}$$

erfüllen  $\mathbb{E}[\hat{a}] = a, \quad \mathbb{E}[\hat{b}] = b.$

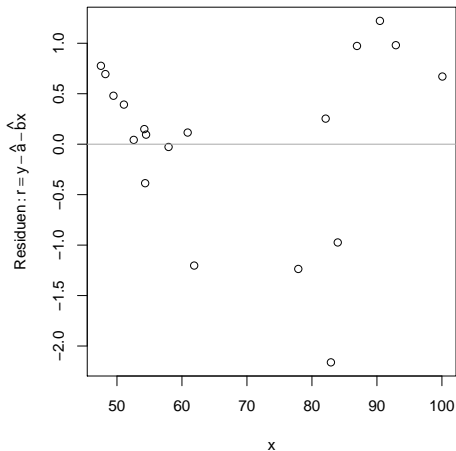
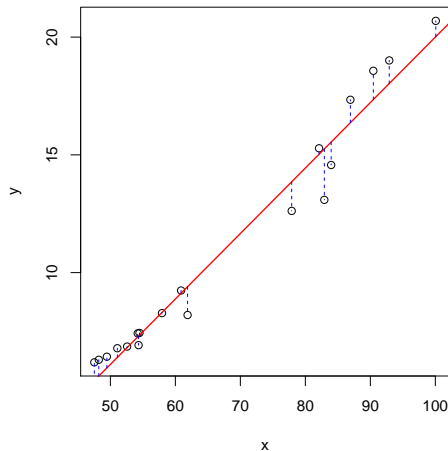
(D.h. sie sind sog. *erwartungstreue Schätzer*.)

Wir schätzen  $\sigma^2$  mit Hilfe der beobachteten Residuenvarianz durch

$$s_{\text{res}}^2 := \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{a} - \hat{b} \cdot x_i)^2. \quad \text{Es gilt } \mathbb{E}[s_{\text{res}}^2] = \sigma^2.$$

(Beachte, dass durch  $n - 2$  geteilt wird. Das hat damit zu tun, dass zwei Modellparameter  $a$  und  $b$  bereits geschätzt wurden, und somit 2 Freiheitsgrade verloren gegangen sind.)

## Geier-Beispiel: Regressionsgerade, Residuen gegen x-Werte



Wir hatten  $\hat{a} = -7.8$ ,  $\hat{b} = 0.278$ , man findet

$$s_{\text{res}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{a} - \hat{b} \cdot x_i)^2} = 0.91.$$

Wir haben  $\hat{a} = -7.8$ ,  $\hat{b} = 0.278$ ,  $s_{\text{res}} = 0.91$  gefunden.

Mit diesen Informationen können wir die Genauigkeit unserer Schätzung beurteilen

und auch einschätzen, wie genau wir einen *neuen* Beobachtungswert vorhersagen könnten.

Beispiel: Angenommen, bei einer weiteren Messung wurde bei Herzfrequenz  $x = 76$  [1/min]

eine Stoffwechselrate von  $y = 14.3$  [J/(g·h)] gemessen.

Die Vorhersage der Regressionsgerade wäre  $-7.8 + 0.278 \cdot 76 = 13.33$ , d.h. sie weicht um  $14.3 - 13.33 = 0.97$  von der Messung ab.

Ist das ein Grund, an unserem Modell zu zweifeln?

Nein: Wenn  $\sigma \approx s_{\text{res}} = 0.91$  gilt, so beobachten wir ein  $\varepsilon_{n+1}$ , das von derselben Größenordnung wie seine Streuung ist — was im Modell mit  $W$ 'keit ca.  $1/3$  passieren kann.

(Wenn wir dagegen  $y = 16.3$  gemessen hätten, wären wir schon beunruhigt ...)

## Beispiel: Rothirsch (*Cervus elaphus*)



photo (c) BS Thurner Hof

Hängt der Anteil männlicher Nachkommen einer Hirschkuh mit ihrem sozialen Rang zusammen?

Frage: Hängt der Anteil männlicher Nachkommen einer Hirschkuh mit ihrem sozialen Rang zusammen?

Betrachten wir folgende Theorie:

Hirschkühe können das Geschlecht ihrer Nachkommen beeinflussen.

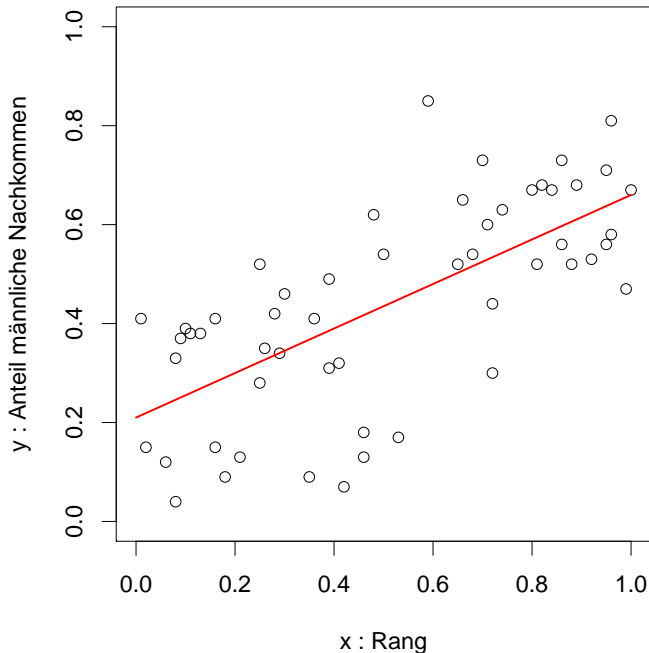
Unter dem Gesichtspunkt evolutionär stabiler Strategien ist zu erwarten, dass schwache Tiere eher zu weiblichem und starke Tiere eher zu männlichem Nachwuchs tendieren.

Folgender Artikel berichtet über eine Langzeitstudie, bei der eine Gruppe Rothirsche auf der schottischen Insel Rùm über 15 Jahre beobachtet wurde, und die zu obiger Frage Daten gesammelt hat:



Clutton-Brock, T. H. , Albon, S. D., Guinness, F. E. (1986)  
Great expectations: dominance, breeding success and  
offspring sex ratios in red deer.  
*Anim. Behav.* **34**, 460—471.

	Rang	Anteil männl. Nachkommen	
1	0.01	0.41	Die Beobachtungen: Für 54 Hirschkühe wurde der Rang (normiert auf einen Wert in $[0, 1]$ , grob gesprochen ein Schätzwert für die Wahrscheinlichkeit, dass die betreffende Kuh ein „Duell“ mit einer zufällig ausgewählten anderen Kuh „gewinnt“) und der Anteil männlicher Nachkommen beobachtet.
2	0.02	0.15	
3	0.06	0.12	
4	0.08	0.04	
5	0.08	0.33	
6	0.09	0.37	
·	·	·	
·	·	·	
·	·	·	
52	0.96	0.81	(Simulierte Daten, die sich an den Werten aus der Originalpublikation orientieren.)
53	0.99	0.47	
54	1.00	0.67	



Es ist

$$\bar{x} = 0.51$$

(mittl. Rang)

$$\bar{y} = 0.44$$

(mittl. Ant. männl. Nachk.)

$$\sigma_x^2 = 0.097$$

$$\text{cov}(x, y) = 0.044$$

$$\hat{b} = \frac{0.044}{0.097} \doteq 0.45$$

$$\hat{a}$$

$$= 0.44 - 0.51 \cdot 0.45$$

$$\doteq 0.21$$



Wir beobachten einen wachsenden (und ungefähr linearen) Zusammenhang zwischen Rang und Anteil männlicher Nachkommen einer Hirschkuh. Ist das ein systematischer Effekt oder könnte reiner Zufall diese Beobachtung genauso gut erklären?

Dazu betrachten wir unser Modell:

$$Y = a + b \cdot X + \varepsilon \quad \text{mit } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Wie berechnet man die Signifikanz eines Zusammenhangs zwischen dem *erklärenden Merkmal*  $X$  und der *Zielgröße*  $Y$ ?

Wir haben  $b$  durch  $\hat{b} = 0.45 \neq 0$  geschätzt. Könnte das wahre  $b$  auch 0 sein?

Anders formuliert: Mit welchem Test können wir der Nullhypothese  $b = 0$  zu Leibe rücken?

Wie groß ist der Standardfehler unserer Schätzung  $\hat{b}$ ?

Modell:

$$y_i = a + b \cdot x_i + \varepsilon_i \quad \text{mit } \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

nicht zufällig:  $a, b, x_i, \sigma^2$       zufällig:  $\varepsilon_i, y_i$

$$\text{var}(y_i) = \text{var}(a + b \cdot x_i + \varepsilon_i) = \text{var}(\varepsilon_i) = \sigma^2$$

und  $y_1, y_2, \dots, y_n$  sind stochastisch unabhängig.

$$\hat{b} = \frac{\sum_i y_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}$$

Demnach

$$\begin{aligned} \text{Var}(\hat{b}) &= \text{Var}\left(\frac{\sum_i y_i (x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}\right) = \frac{\text{Var}(\sum_i y_i (x_i - \bar{x}))}{(\sum_i (x_i - \bar{x})^2)^2} \\ &= \frac{\sum_i \text{Var}(y_i) (x_i - \bar{x})^2}{(\sum_i (x_i - \bar{x})^2)^2} = \sigma^2 \cdot \frac{\sum_i (x_i - \bar{x})^2}{(\sum_i (x_i - \bar{x})^2)^2} \\ &= \sigma^2 / \sum_i (x_i - \bar{x})^2 \end{aligned}$$

Tatsächlich ist  $\hat{b}$  Normalverteilt mit Mittelwert  $b$  und

$$\text{Var}(\hat{b}) = \sigma^2 / \sum_i (x_i - \bar{x})^2$$

**Problem:** Wir kennen  $\sigma^2$  nicht.

Wir schätzen  $\sigma^2$  mit Hilfe der beobachteten Residuenvarianz durch

$$s_{\text{res}}^2 := \frac{\sum_i (y_i - \hat{a} - \hat{b} \cdot x_i)^2}{n - 2}$$

Erinnerung: Hier wird durch  $n - 2$  geteilt. Das hat damit zu tun, dass zwei Modellparameter  $a$  und  $b$  bereits geschätzt wurden, und somit 2 Freiheitsgrade verloren gegangen sind.

$$\text{Var}(\hat{b}) = \sigma^2 / \sum_i (x_i - \bar{x})^2$$

Schätze  $\sigma^2$  durch

$$s_{\text{res}}^2 = \frac{\sum_i (y_i - \hat{a} - \hat{b} \cdot x_i)^2}{n - 2}.$$

Dann ist der Standardfehler von  $\hat{b}$  gegeben durch  $\frac{s_{\text{res}}}{\sqrt{\sum_i (x_i - \bar{x})^2}}$   
 (unser Schätzwert für die Streuung von  $\hat{b}$ )  
 und (unter den Modellannahmen)

$$\frac{\hat{b} - b}{s_{\text{res}} / \sqrt{\sum_i (x_i - \bar{x})^2}} = \frac{\hat{b} - b}{s_{\text{res}} / \sigma_x \sqrt{n}}$$

Student- $t$ -verteilt mit  $n - 2$  Freiheitsgraden. Wir können also den  $t$ -Test anwenden, um die Nullhypothese  $b = 0$  zu testen.

Im Rothirschkühe-Beispiel:

$$\hat{b} = 0.45, \frac{s_{\text{res}}}{\sqrt{\sum_i (x_i - \bar{x})^2}} = 0.0673,$$

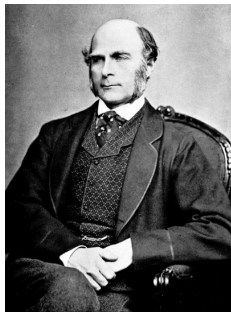
also beobachten wir  $t = \frac{\hat{b} - 0}{s_{\text{res}} / \sqrt{\sum_i (x_i - \bar{x})^2}} = 6.7$

Einer Tabelle entnehmen wir: Das 99.95%-Quantil der Student-Verteilung mit 50 Freiheitsgraden ist 3.496 (und das der Student-Vert. mit 60 Freiheitsgraden ist 3.460).

Wir können also die Nullhypothese „das wahre  $b = 0$ “ zum Signifikanzniveau 0.1% ablehnen.

Bemerkung: Das beweist natürlich nicht, dass Hirschkühe das Geschlecht ihrer Nachkommen willentlich bestimmen können. Es scheint eher plausibel anzunehmen, dass es Faktoren gibt, die den Rang und die Geschlechterverteilung der Nachkommen zugleich beeinflussen, siehe die Diskussion in dem zitierten Artikel von T. H. Clutton-Brock et. al.

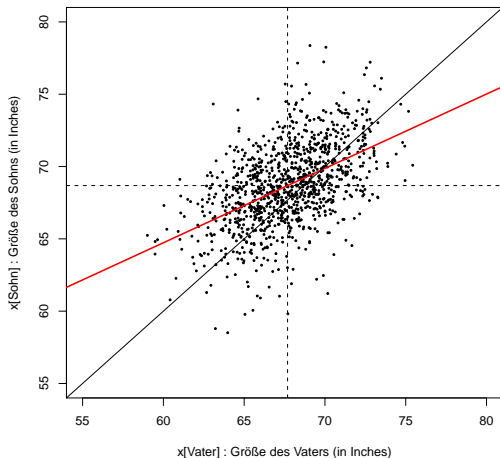
Woher kommt der Name „Regression“ (nach lat. regressio, Zurückkommen)?



Francis Galton (1822–1911, engl. Wissenschaftler) hat angesichts biometrischer Beobachtungen den Ausdruck “regression towards the mean” geprägt.

## Ein (relativ berühmter) Datensatz von Karl Pearson (1858-1936)

1078 Größen von Vater und Sohn

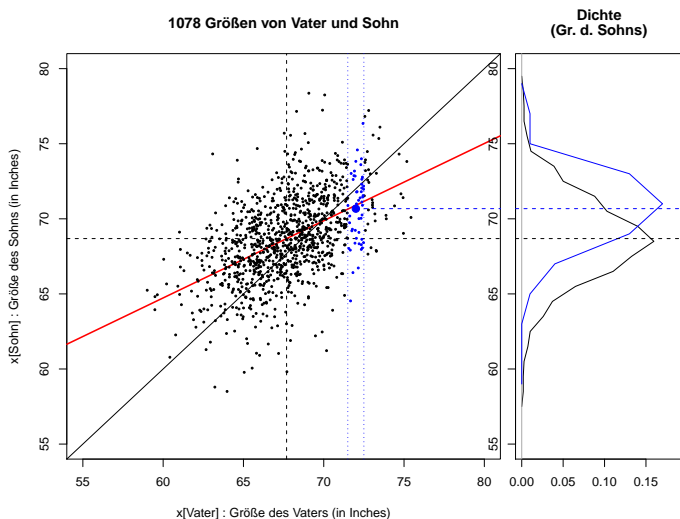


$$\bar{x}_{\text{Vater}} = 67.7, \bar{x}_{\text{Sohn}} = 68.7, \sigma_{\text{Vater}}^2 = 7.52 \quad (\sigma_{\text{Vater}} = 2.74, \sigma_{\text{Sohn}} = 2.81),$$

$$\text{COV}(x_{\text{Vater}}, x_{\text{Sohn}}) = 3.87$$

$$(\text{Korrelationskoeffizient } \rho = \text{COV}(x_{\text{Vater}}, x_{\text{Sohn}}) / (\sigma_{\text{Vater}} \sigma_{\text{Sohn}})) = 0.50)$$

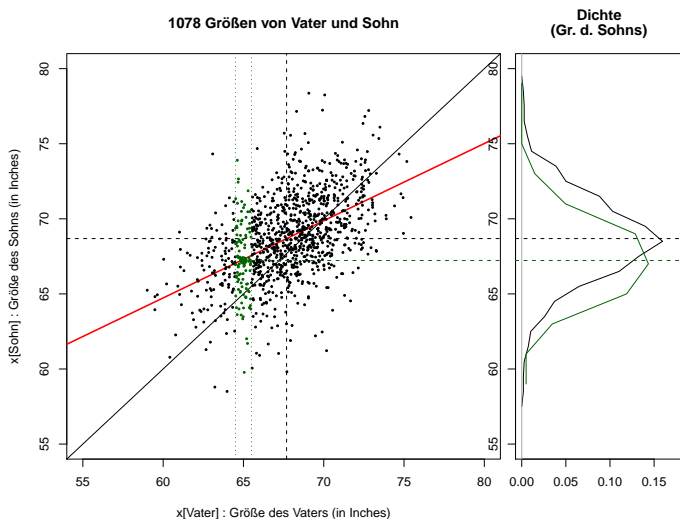
$$\text{Regressionsgerade: } x_{\text{Sohn}} = 33.89 + 0.514x_{\text{Vater}}.$$



Betrachten wir die Söhne von überdurchschnittlich großen Vätern (z.B. Väter, die ca. 72 Inches groß sind):

Diese Söhne sind überdurchschnittlich groß (im Vergleich zu allen Söhnen), aber im Mittel kleiner als ihr Vater.

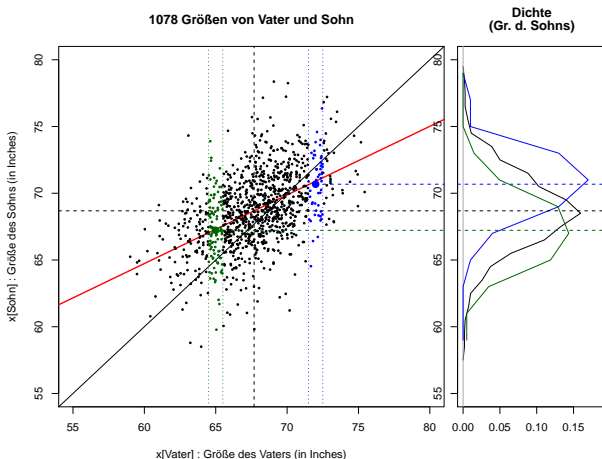




Betrachten wir andererseits die Söhne von unterdurchschnittlich großen Vätern (z.B. Väter, die ca. 65 Inches groß sind):

Diese Söhne sind unterdurchschnittlich groß (im Vergleich zu allen Söhnen), aber im Mittel größer als ihr Vater.

# “Regression towards the mean”



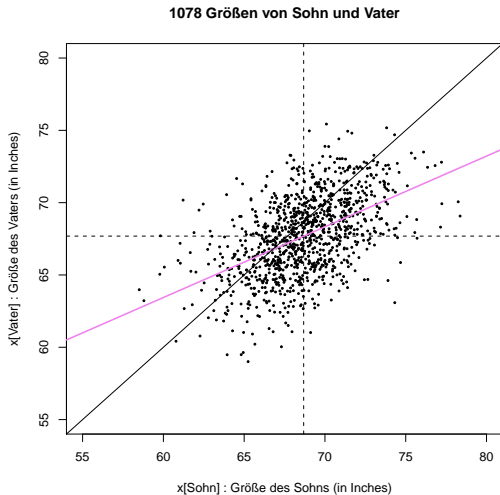
Wir sehen: Söhne überdurchschnittlich großer Väter sind im Mittel kleiner als ihr Vater (aber immer noch größer als der Populationsdurchschnitt), für Söhne unterdurchschnittlich großer Väter ist es umgekehrt: „Rückkehr zum Mittelwert“.

Bemerkung: Das beobachtete Phänomen der „Rückkehr zum Mittelwert“ bedeutet nicht notwendigerweise einen tieferen kausalen Zusammenhang, es tritt stets im Zusammenhang mit natürlicher Variabilität auf (technisch gesehen stets, wenn für den Korrelationskoeffizient  $\rho$  gilt  $|\rho| < 1$ ).

Bestimmen wir (spañeshalber) im Größen-Beispiel die Regressionsgerade für die Größe des Vaters als Funktion der Größe des Sohns:

Wir hatten  $\bar{x}_{\text{Vater}} = 67.7$ ,  $\bar{x}_{\text{Sohn}} = 68.7$ ,  $\text{cov}(x_{\text{Vater}}, x_{\text{Sohn}}) = 3.87$ ,  
 $\sigma_{\text{Sohn}}^2 = 7.92$

und finden die Regressionsgerade  $x_{\text{Vater}} = 34.1 + 0.489x_{\text{Sohn}}$

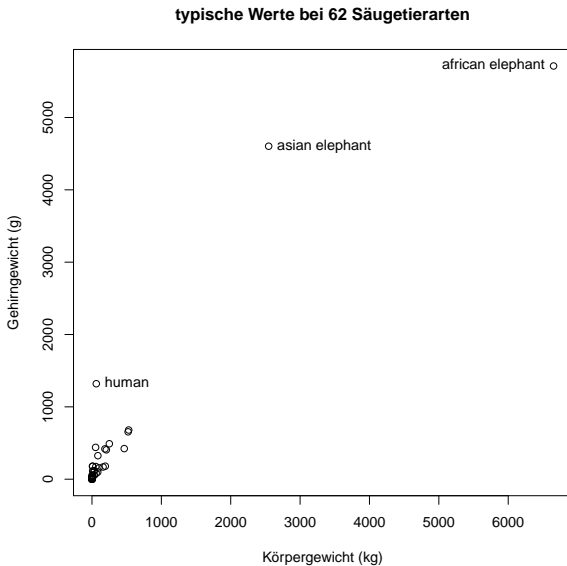


Regressionsgerade:  $x_{\text{Vater}} = 34.1 + 0.489x_{\text{Sohn}}$

## Daten: Typisches Körpergewicht [kg] und Gehirngewicht [g] von 62 Säugetierarten

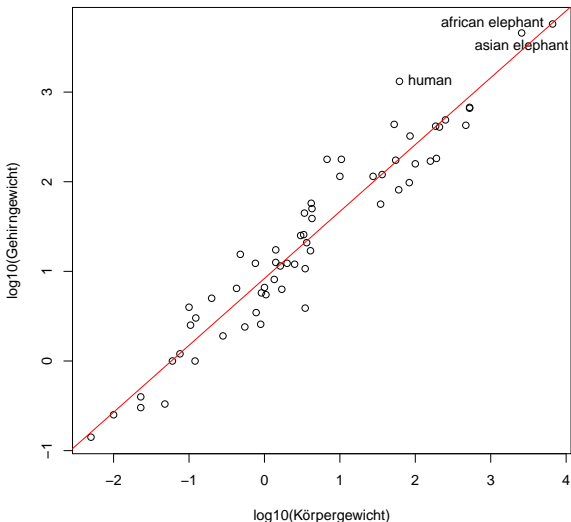
	weight.kg.	brain.weight.g	species
1	6654.00	5712.00	african elephant
2	1.00	6.60	
3	3.39	44.50	
4	0.92	5.70	
5	2547.00	4603.00	asian elephant
6	10.55	179.50	
7	0.02	0.30	
8	160.00	169.00	
9	3.30	25.60	cat
10	52.16	440.00	chimpanzee
11	0.43	6.40	
.	.	.	.
.	.	.	.

(Daten nach H.J. Jerison, Evolution of the brain and intelligence. Academic Press, New York, 1973.)



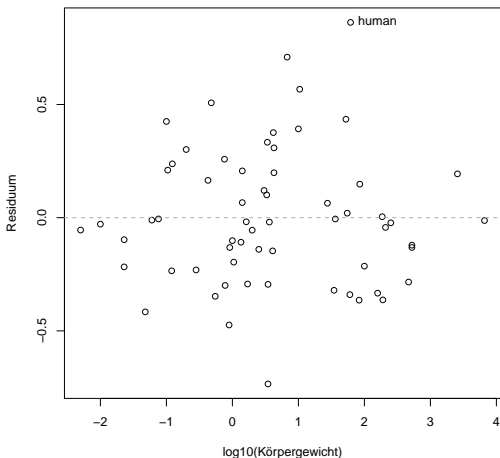
Wir sehen — wenig: Wenige „schwere“ Arten dominieren die Skala.  
Abhilfe: Logarithmieren!

## typische Werte bei 62 Säugetierarten



Auf der log-Skala sieht ein linearer Zusammenhang plausibel aus.

Die Regressionsgerade ist  $y = 0.921 + 0.746x$ .



Die Residuen passen in etwa zur Modellvorstellung  
„ $\log(\text{Gehirngewicht}) = 0.921 + 0.746 \cdot \log(\text{Körpergewicht}) + \varepsilon$ “,  
wobei die Streuung des „Fehlerterms“  $\varepsilon$  nicht von der Art abhängt.



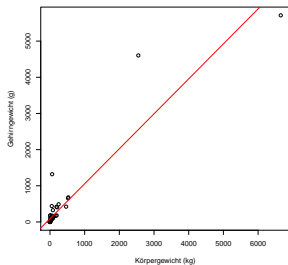
Wir haben zunächst die  $x$ -Werte (Körpergewichte) und die  $y$ -Werte (Gehirngewichte) logarithmiert, und dann die Regressionsgerade angepasst.

Frage: War das nur ein Trick, damit das Diagramm besser erkennbar erkennbar wird, oder gibt es eine weitere inhaltliche Bedeutung?

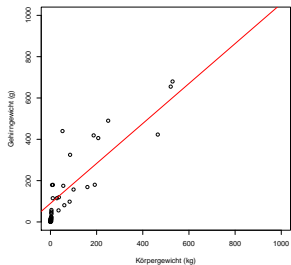
Passen wir den unlogarithmierten Wertepaaren ebenfalls eine Regressionsgerade an

(es kommt  $y_{\text{Gehirngew.}} = 89.91 + 0.967x_{\text{Körpergew.}}$  heraus).

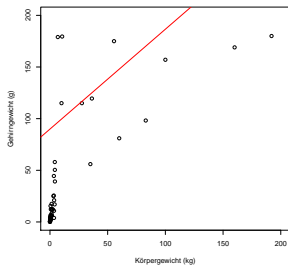
typische Werte bei 62 Säugetierarten



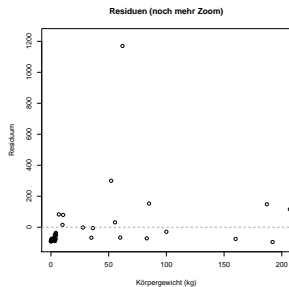
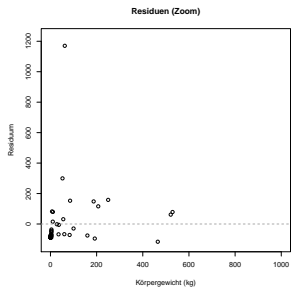
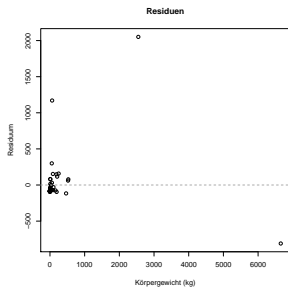
typische Werte bei 62 Säugetierarten (Zoom)



typische Werte bei 62 Säugetierarten (noch mehr Zoom)



Die Anpassung wirkt schlechter. Wir betrachten die „leichteren“ Arten näher. Und noch näher.



Die Residuen wirken nun sehr inhomogen, was nicht zur Modellvorstellung passt. (Der Eindruck bleibt auch bei näherer Betrachtung der leichten Arten bestehen und auch bei noch näherer.)

Wir sehen, dass die Varianz der Residuen von den angepassten Werten bzw. dem Körpergewicht abhängt. Man sagt, es liegt *Heteroskedastizität* vor.

Das Modell geht aber von *Homoskedastizität* aus, d.h. die Residuenvarianz soll von den erklärenden Merkmalen (dem Körpergewicht) und den angepassten Werten (annähernd) unabhängig sein.

### **Varianzstabilisierende Transformation:**

Wie können wir die Körper- und Hirnmasse umskalieren, um Homoskedastizität zu erreichen?

Eigentlich ist es ja offensichtlich: Bei Elefanten kann das typischerweise 5 kg schwere Hirn je nach Individuum auch mal 500 g schwerer oder leichter sein. Wenn bei einer Tierart das Hirn typischerweise 5 g schwer ist, wird es nicht um 500 g variieren können, sondern vielleicht ebenfalls um 10%, also  $\pm 0.5$  g. Die Variabilität ist hier also nicht additiv, sondern multiplikativ:

$$\text{Hirnmasse} = (\text{erwartete Hirnmasse}) \cdot \text{„Rauschen“}$$

Das können wir aber in etwas mit additivem Zufallsterm umwandeln, indem wir auf beiden Seiten den Logarithmus anwenden:

$$\log(\text{Hirnmasse}) = \log(\text{erwartete Hirnmasse}) + \log(\text{„Rauschen“})$$