Biostatistik

Lineare Regression

Matthias Birkner



Inhalt

- Lineare Regression: wozu und wie?
- t-Test für lineare Zusammenhänge
- Weitere Beispiele und Anmerkungen
 - K. Pearsons Größen von Vätern und Söhnen
 - Körper- und Gehirngewicht: Transformation der Daten
- Lineare Regression mit R



Gypus fulvus Gänsegeier

photo (c) by Jörg Hempel

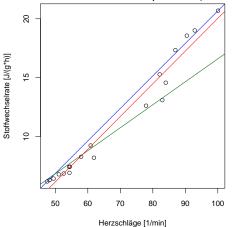
Prinzinger, R., E. Karl, R. Bögel, Ch. Walzer (1999): Energy metabolism, body temperature, and cardiac work in the Griffon vulture Gyps vulvus - telemetric investigations in the laboratory and in the field.

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:

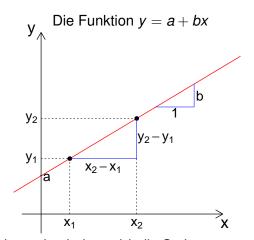
19 Beobachtungen eines Geiers im Labor bei konstanter Temperatur (16 $^{\circ}$ C)



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

Welche Gerade passt "am Besten"?

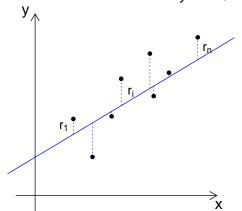
Erinnerung: Lineare Funktionen



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 $\frac{y_2 - y_1}{x_2 - y_1} = b$.

Unsere Situation: n beobachtete Paare (x_i, y_i) , i = 1, ..., n

Die Daten und eine Gerade y = a + bx

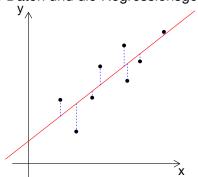


Für keine gegebene Steigung b und Achsenabschnitt a liegen die Beobachtungen genau auf der Geraden y=a+bx.

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

kleinste-Quadrate-Schätzer \widehat{a} , \widehat{b} und Regressionsgerade

Die Daten und die Regressionsgerade



Ansatz: Finde \widehat{a} und \widehat{b} derart, dass $\sum_{i=1}^n \left(y_i - \widehat{a} - \widehat{b}\,x_i\right)^2 \stackrel{!}{=}$ 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 $\overline{x} := \frac{1}{n} \sum_{i=1}^{n} x_i$, $\overline{y} := \frac{1}{n} \sum_{i=1}^{n} y_i$ die (empirischen) Mittelwerte der x- bzw. der y-Werte,

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

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

Satz (Koeffizienten der Regressionsgerade)

$$\widehat{b} = \frac{\operatorname{cov}(x,y)}{\sigma^2}, \quad \widehat{a} = \overline{y} - \widehat{b}\,\overline{x}.$$

Bemerkungen zu den Regressionskoeffizienten

$$\widehat{b} = \frac{\operatorname{cov}(x, y)}{\sigma_x^2}, \quad \widehat{a} = \overline{y} - \widehat{b}\overline{x}$$

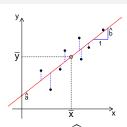
- Beobachtung: Die Regressionsgerade $y = \hat{a} + \hat{b}x$ geht durch den Schwerpunkt der Daten (\bar{x}, \bar{y}) . (Dies kann eine Merkhilfe für die Formeln bilden.)
- 2 Man kann \hat{b} auch anders ausdrücken:

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

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

Bemerkungen zu den Regressionskoeffizienten

$$\widehat{b} = \frac{\operatorname{cov}(x, y)}{\sigma^2}, \quad \widehat{a} = \overline{y} - \widehat{b}\overline{x}$$



- 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:

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

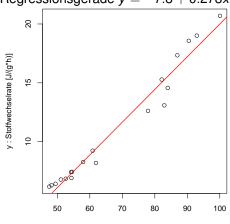
(Für das 1. Gleichheitszeichen erweitere mit n, für das 2. beachte $\overline{y} \sum_{i=1}^{n} (x_i - \overline{x}) = 0$. Die varianten Formeln können ggf. zum Rechnen angenehmer sein.)

Für das Geier-Beispiel ist mit x = Herzfrequenz, y = Stoffwechselrate:

$$\overline{x} = 67.9, \overline{y} = 11.1, cov(x, y) = 83.5, \sigma_x^2 = 300.7, also$$

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

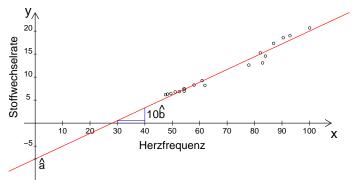
n = 19 Datenpunkte und Regressionsgerade y = -7.8 + 0.278x



x: Herzschläge [1/min]

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?

$$\frac{1}{n}\sum_{i=1}^{n}r_i^2 = \frac{1}{n}\sum_{i=1}^{n}\left(y_i - a - bx_i\right)^2$$
$$= \left(\frac{1}{n}\sum_{i=1}^{n}\left(y_i - \overline{y} - b(x_i - \overline{x})\right)^2\right) + \left(\overline{y} - a - b\overline{x}\right)^2$$

denn

$$(y_i - \overline{y} - b(x_i - \overline{x}))^2 = (y_i - a - bx_i - (\overline{y} - a - b\overline{x}))^2$$

= $(y_i - a - bx_i)^2 - 2(y_i - a - bx_i)(\overline{y} - a - b\overline{x}) + (\overline{y} - a - b\overline{x})^2$

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

Regressionskoeffizienten: Woher kommen die Formeln?

$$\frac{1}{n} \sum_{i=1}^{n} r_{i}^{2} = \frac{1}{n} \sum_{i=1}^{n} (y_{i} - a - bx_{i})^{2}$$

$$= \frac{1}{n} \sum_{i=1}^{n} (y_{i} - \overline{y} - b(x_{i} - \overline{x}))^{2} + (\overline{y} - a - b\overline{x})^{2}$$

$$= \frac{1}{n} \sum_{i=1}^{n} (y_{i} - \overline{y})^{2} - 2b \frac{1}{n} \sum_{i=1}^{n} (x_{i} - \overline{x})(y_{i} - \overline{y})$$

$$+ b^{2} \frac{1}{n} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2} + (\overline{y} - a - b\overline{x})^{2}$$

$$= \sigma_{y}^{2} - 2b \operatorname{cov}(x, y) + b^{2} \sigma_{x}^{2} + (\overline{y} - a - b\overline{x})^{2}$$

 $= \sigma_y^2 - \frac{(\operatorname{cov}(x,y))^2}{\sigma_x^2} + \sigma_x^2 \Big(b - \frac{\operatorname{cov}(x,y)}{\sigma_x^2} \Big)^2 + (\overline{y} - a - b\overline{x})^2$

Regressionskoeffizienten: Woher kommen die Formeln? Demnach:

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

$$= \underbrace{\sigma_y^2 - \frac{(\text{Cov}(x, y))^2}{\sigma_x^2}}_{\text{h\tilde{a}ngt nicht von } a, b} + \sigma_x^2 \underbrace{\left(b - \frac{\text{Cov}(x, y)}{\sigma_x^2}\right)^2}_{\geq 0} + \underbrace{\left(\overline{y} - a - b\overline{x}\right)^2}_{\geq 0}$$

Die Summe der Residuenquadrate wird also minimiert durch die Wahlen

$$\widehat{b} = \frac{\operatorname{Cov}(x, y)}{\sigma^2}, \quad \widehat{a} = \overline{y} - \widehat{b}\,\overline{x}.$$

(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 ε_i unabhängig und normalverteilt sind mit Mittelwert 0 und derselben Varianz σ^2 .

gegebene Daten:

 $\begin{array}{cccc}
\mathbf{Y} & \mathbf{X} \\
y_1 & x_1 \\
y_2 & x_2 \\
y_3 & x_3 \\
\vdots & \vdots
\end{array}$

 y_n

 X_n

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

$$y_{1} = a + b \cdot x_{1} + \varepsilon_{1}$$

$$y_{2} = a + b \cdot x_{2} + \varepsilon_{2}$$

$$y_{3} = a + b \cdot x_{3} + \varepsilon_{3}$$

$$\vdots$$

$$\vdots$$

$$y_{n} = a + b \cdot x_{n} + \varepsilon_{n}$$

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, σ^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, σ^2 zufällig: ε_i, y_i

Die kleinste-Quadrate-Schätzer

$$\hat{b} = \frac{\sum_{i} y_{i}(x_{i} - \overline{x})}{\sum_{i} (x_{i} - \overline{x})^{2}}, \quad \hat{a} = \overline{y} - b\overline{x}$$

erfüllen

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

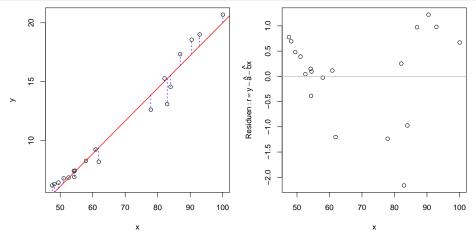
(D.h. sie sind sog. erwartungstreue Schätzer.)

Wir schätzen σ^2 mit Hilfe der beobachten Residuenvarianz durch

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

(Beachte, dass durch n-2 geteilt wird. Das hat damit zu tun, dass zwei Modellparameter a und b bereit 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

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

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

eine Stoffwechselrate von y=14.3 [J/(g·h)] gemessen. Die Vorhersage der Regressionsgerade wäre

Herzfrequenz x = 76 [1/min]

 $-7.8 + 0.278 \cdot 76 = 13.33$, 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 ε_{n+1} , das von derselben Größenordnung wie seine Streuung ist — was im Modell mit W'keit en 1/2 passioren kann

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.

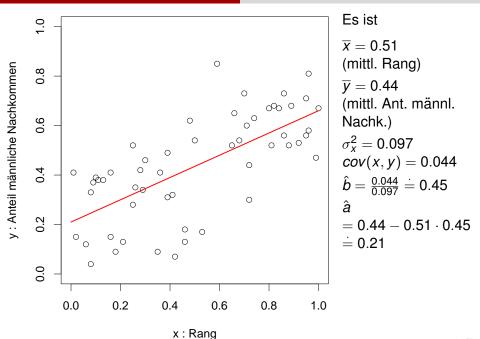
Unter dem Gesichtspunkt evolutionär stabiler Strategien ist zu erwarten, dass schwache Tiere eher zu weiblichem und starke Tiere eher zu männlichem Nachwuchs tendieren.

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.

1 2 3 4 5 6	Rang 0.01 0.02 0.06 0.08 0.08 0.09	Anteil männl. 0.41 0.15 0.12 0.04 0.33 0.37 .	Nachkommen 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 Nach-			
•			kommen beobachtet.			
53	0.96 0.99 1.00	0.81 0.47 0.67	(Simulierte Daten, die sich an den Werten aus der Originalpublikation orientieren.)			



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 genausogut erklären?

Dazu betrachten wir unser Modell:

$$Y = a + b \cdot X + \varepsilon$$
 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:

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

nicht zufällig: a, b, x_i , σ^2 zufällig: ε_i , y_i

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

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

$$Var(\hat{b}) = Var\left(\frac{\sum_{i} y_{i}(x_{i} - \bar{x})}{\sum_{i} (x_{i} - \bar{x})^{2}}\right) = \frac{Var\left(\sum_{i} y_{i}(x_{i} - \bar{x})\right)}{\left(\sum_{i} (x_{i} - \bar{x})^{2}\right)^{2}}$$

$$= \frac{\sum_{i} Var(y_{i}) (x_{i} - \bar{x})^{2}}{\left(\sum_{i} (x_{i} - \bar{x})^{2}\right)^{2}} = \sigma^{2} \cdot \frac{\sum_{i} (x_{i} - \bar{x})^{2}}{\left(\sum_{i} (x_{i} - \bar{x})^{2}\right)^{2}}$$

$$= \sigma^{2} / \sum_{i} (x_{i} - \bar{x})^{2}$$

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

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

Problem: Wir kennen σ^2 nicht.

Wir schätzen σ^2 mit Hilfe der beobachten Residuenvarianz durch

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

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

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

Schätze σ^2 durch

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

Dann ist der Standardfehler von \hat{b} gegeben durch $\frac{s_{\rm res}}{\sqrt{\sum_i (x_i - \bar{x})^2}}$ (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:

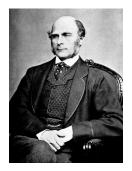
$$\hat{b} = 0.45, \, rac{s_{
m res}}{\sqrt{\sum_i (x_i - ar{x})^2}} = 0.0673,$$

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.

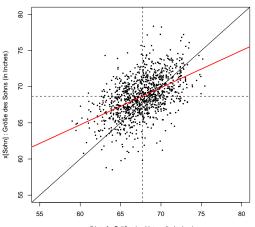
Bemerkung: Das beweist natürlich nicht, dass Hirschkühe das Geschlecht ihrer Nachkommen willentlich bestimmen können. Es scheint eher plausibel anzunehmen, dass es Faktoren gibt, die den Rang und die Geschlechterverteilung der Nachkommen zugleich beeinflussen, siehe die Diskussion in dem zitierten Artikel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

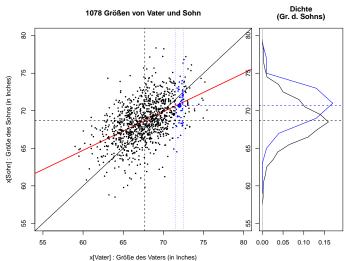
^{*}Siehe auch Robert Langkjaer-Bain, The troubling legacy of Francis Galton, Significance 16(3), 2019.

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



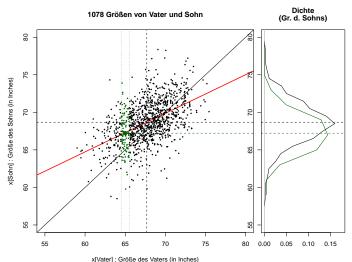
$$\overline{x}_{ ext{Vater}} = 67.7, \, \overline{x}_{ ext{Sohn}} = 68.7, \, \sigma_{ ext{Vater}}^2 = 7.52 \, (\sigma_{ ext{Vater}} = 2.74, \, \sigma_{ ext{Sohn}} = 2.81), \, \cos(x_{ ext{Vater}}, x_{ ext{Sohn}}) = 3.87 \, (\text{Korrelationskoeffizient } \rho = \cos(x_{ ext{Vater}}, x_{ ext{Sohn}})/(\sigma_{ ext{Vater}}, \sigma_{ ext{Sohn}})) = 0.50)$$

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



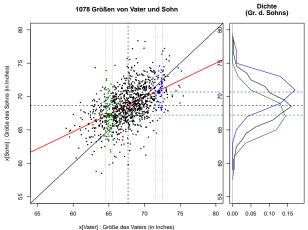
Betrachten wir die Söhne von überdurchschnittlich großen Vätern (z.B. Väter, die ca. 72 Inches groß sind):

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".

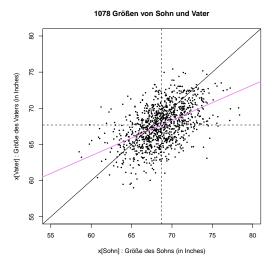
38/52

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 Korrelationkoeffizient ρ gilt $|\rho| < 1$).

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

Wir hatten $\overline{x}_{Vater} = 67.7$, $\overline{x}_{Sohn} = 68.7$, $cov(x_{Vater}, x_{Sohn}) = 3.87$, $\sigma_{Sohn}^2 = 7.92$ und finden die Regressionsgerade $x_{Vater} = 34.1 + 0.489x_{Sohn}$

39/52



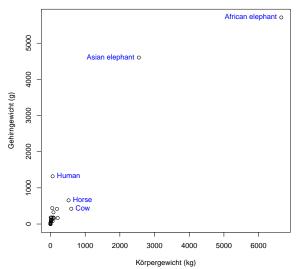
Regressionsgerade: $x_{\text{Vater}} = 34.1 + 0.489 x_{\text{Sohn}}$

Daten: Typisches Körpergewicht [kg] und Gehirngewicht [g] von 56 Säugetierarten

9					
	brain.weight.g			brain.weight.g	
Short-nosed echidna	25.00	4.500	Cotton rat	1.18	0.148
North American Opossum	6.30	1.700	Golden hamster	1.00	0.120
Phalanger	11.40	1.620	Mole rat	3.00	0.122
European hedgehog	3.50	0.770	African giant pouched rat	6.60	1.000
Desert hedgehog	2.40	0.550	Laboratory rat	1.90	0.320
Tenrec	2.60	0.900	House mouse	0.40	0.022
Greater short-tailed shrew	0.29	0.019	Guinea pig	5.50	0.728
Lesser short-tailed shrew	0.14	0.005	Chinchilla	6.40	0.420
Musk shrew	0.33	0.048	Red fox	50.40	4.230
Star-nosed mole	1.00	0.060	Arctic fox	44.50	3.380
Eastern american mole	1.20	0.075	Dog	70.00	14.000
Tree shrew	2.50	0.104	Genet	17.50	2.000
Little brown bat	0.25	0.010	Domestic cat	25.60	3.300
Big brown bat	0.30	0.023	Jaguar	157.00	100.000
Slow loris	12.50	1.400	Gray seal	325.00	85.000
Galago	5.00	0.200	Asian elephant	4603.00	2547.000
Squirrel monkey	20.00	0.743	African elephant	5712.00	6654.000
Owl monkey	15.50	0.480	Rock hyrax	21.00	3.600
Patas monkey	115.00	10.000	Gray hyrax	12.27	2.625
Macaque	179.00	6.800	Tree hyrax	12.30	2.950
Baboon	180.00	25.235	Horse	655.00	521.000
Chimpanzee	440.00	52.200	Donkey	419.00	187.000
Human	1320.00	62.000	Brazilian tapir	169.00	207.501
Giant armadillo	81.00	60.000	Pig	180.00	86.250
Long-nosed armadillo	10.80	3.500	Roe deer	98.20	14.800
Rabbit	12.10	2.500	Cow	423.00	600.000
Arctic ground squirrel	5.70	0.920	Goat	115.00	33.500
Thirteen-lined ground squirrel	4.00	0.101	Sheep	175.00	55.500

Daten entnommen aus V. M. Savage and G. B. West. A quantitative, theoretical framework for understanding mammalian sleep. Proceedings of the National Academy of Sciences, 104 (3):1051-1056, 2007 (nur Spezies, für die Information zum mittleren Gehirngewicht vorlag)

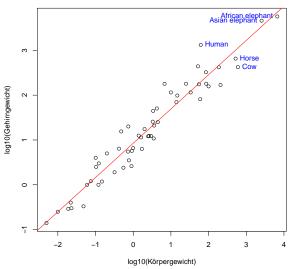
typische Werte bei 56 Säugetierarten



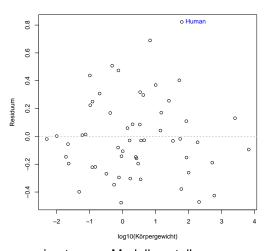
Wir sehen — wenig: Wenige "schwere" Arten dominieren die Skala.

Abhilfe: Logarithmieren!

typische Werte bei 56 Säugetierarten



Auf der log-Skala sieht ein linearer Zusammenhang plausibel aus. Die Regressionsgerade ist y = 0.926 + 0.765x.



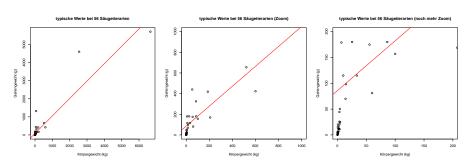
Die Residuen passen in etwa zur Modellvorstellung "log(Gehirngewicht)=0.926+0.765 \cdot log(Körpergewicht) + ε ", wobei die Streuung des "Fehlerterms" ε nicht von der Art abhängt.

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

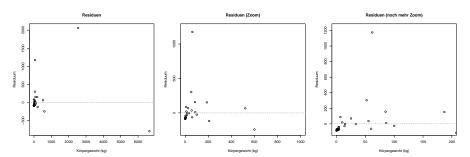
Frage: War das nur ein Trick, damit das Diagramm besser erkennbar erkennbar wird, oder gibt es eine weitere inhaltliche Bedeutung?

Passen wir den unlogarithmierten Wertepaaren ebenfalls eine Regressiongerade an

(es kommt $y_{\text{Gehirngew.}} = 85.92 + 0.964 x_{\text{K\"{o}rpergew.}}$ heraus).



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?

Idee: 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 ± 0.5 g. Die Variabilität ist hier also nicht additiv, sondern multiplikativ:

 $\label{eq:hirnmasse} \textit{Hirnmasse}) \cdot \textit{``,Rauschen''}$

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

 $\log(\mathsf{Hirnmasse}) = \log(\mathsf{erwartete}\;\mathsf{Hirnmasse}) + \log(\mathsf{"Rauschen"})$

Lineare Regression mit R

```
Eingabe der Datensätze x und y (z.B. "von Hand"):

x \leftarrow c(1.0, 2.2, 2.7, 2.7, 3.5, 5.0)
```

```
y \leftarrow c(4.9, 7.2, 8.8, 8.4, 10.4, 11.2)
```

Ausgabe von Steigung und *y*-Achsenabschnitt der Ausgleichsgeraden, sowie Kovarianz und Korrelationskoeffizient

```
lm(y \sim x)

cov(x, y)

cor(x, y)
```

Grafische Darstellung

```
plot(x, y) # die Datenpunkte malen abline(lm(y \sim x), col='red', lwd=3) # und Ausgleichsgerade hinzufügen
```