

Stochastik für Informatiker

Einfache lineare Regression

Hanno Gottschalk

July 4, 2023

Streudiagramme	3
Bremsweg - 1920	4
Werbung und Verkauf	5
Streudiagramm Werbung und Verkauf	6
Lineare Korrelation	7
Motivation	8
Die emp. Kovarianz.	9
Der Korrelationskoeffizient	10
Eigenschaften des Korrelationskoeffizienten	11
Beispiele für Korrelation	12
Einstufung der Korrelation.	13
Kausalitätsmodelle	14
Kausale Abhängigkeiten	15
Scheinkorrelation	16
Beispiel Scheinkorrelation.	17
Beispiel Scheinkorrelation II	18
Abhängigkeit ohne lin. Korrelation	19
Berechnung von Ausgleichsgraden (lin. Reg.)	20
Vorbemerkung - Modellierung I.	21
Vorbemerkung - Modellierung II	22
Overfitting	23
Lineare Abhängigkeit	24
Kleinste Quadrate (least squares)	25
Lin. Reg. – Berechnung	26
Lin. Reg. – Berechnung	27
Ausgleichsgrade Verkäufe.	28
Effektiv Rechnen	29
Rechenbeispiel	30
Rechnen mit Taschenrechner	31
R – wie Rechnen.	32

Inhaltsverzeichnis Vorlesung

- Streudiagramme
- Lineare Korrelation
- Kausalitätsmodelle
- Berechnung von Ausgleichsgraden (Lin. Reg.)
- Effektiv Rechnen

Hanno Gottschalk

Stochastik für Informatiker – 2 / 32

Streudiagramme

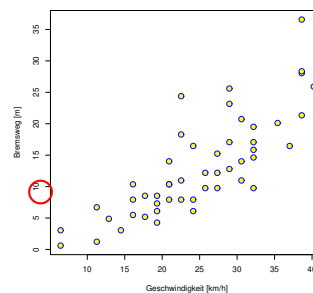
3 / 32

Bremsweg - 1920

- $X : \Omega \rightarrow [0, \infty)$ Geschwindigkeit [kmh], $Y : \Omega \rightarrow [0, \infty)$ Bremsweg[m]
- Ω Bremsvorgänge von PKW 1920
- $\chi \subseteq \Omega$ Bremsvorgänge, die während einer Testkampagne 1920 gemessen wurden.



Bremsweg von Autos in den 1920er Jahren



Je höher die Geschwindigkeit, desto länger der Bremsweg
(aber nicht immer)

Hanno Gottschalk

Stochastik für Informatiker – 4 / 32

Werbung und Verkauf

Monat	1	2	3	4	5	6	7	8	9	10
Werbeausgaben	1.2	0.8	1.0	1.3	0.7	0.8	1.0	0.6	0.9	1.1
Verkäufe	101	92	110	120	90	82	93	75	91	105

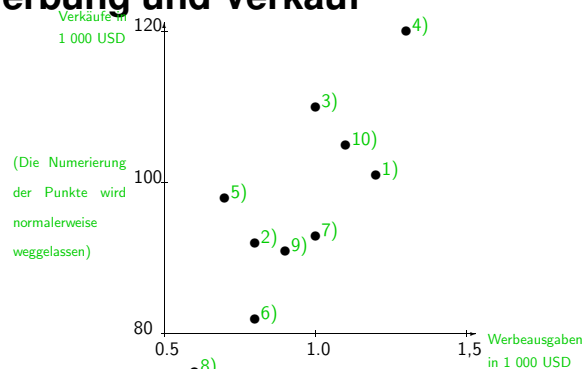
Die Kontingenztabelle (2 Werte fehlen) ist im Wesentlichen leer (und damit weitgehend unbrauchbar)

W.ausg./Verk	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105
1.3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
1.1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
1.0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
0.9	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0.8	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
0.7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0.6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Hanno Gottschalk

Stochastik für Informatiker – 5 / 32

Streudiagramm Werbung und Verkauf



Monat	1	2	3	4	5	6	7	8	9	10
Werbeausgaben	1.2	0.8	1.0	1.3	0.7	0.8	1.0	0.6	0.9	1.1
Verkäufe	101	92	110	120	90	82	93	75	91	105

Besser, da nicht Unmengen redundanter Nullen mitgeschleppt werden ...

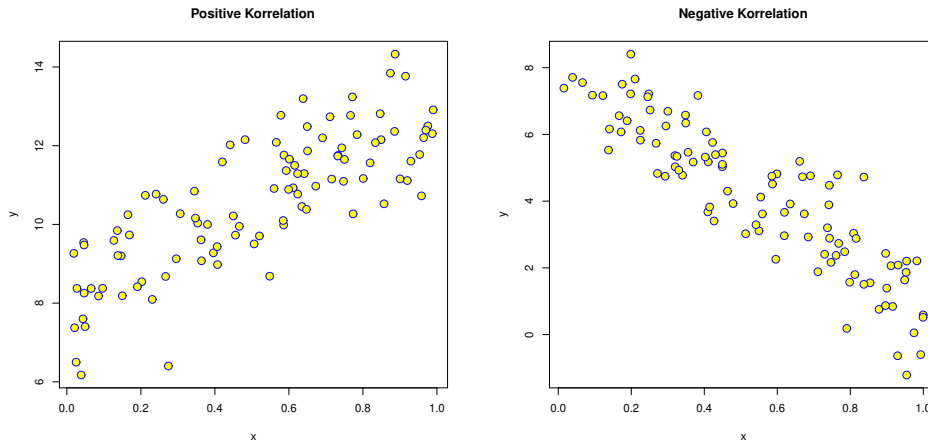
Hanno Gottschalk

Stochastik für Informatiker – 6 / 32

Motivation

Wir sprechen von Korrelation von numerischen Merkmalen X und Y wenn

- Mit Anstieg von X steigt auch Y tendenziell an (*positive Korrelation*)
- Mit Anstieg von X fällt Y tendenziell ab (*negative Korrelation*)



Hanno Gottschalk

Stochastik für Informatiker – 8 / 32

Die emp. Kovarianz

Def.: Es seien in der Stichprobe $\chi \subseteq \Omega$, $\chi = \{\omega_1, \dots, \omega_n\}$ für die Merkmale X und Y die Merkmalsausprägungen x_1, \dots, x_n bzw. y_1, \dots, y_n gemessen worden. Die emp. Kovarianz ist gegeben durch

$$\hat{\sigma}_{X,Y} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{n}{n-1} (\overline{xy} - \bar{x}\bar{y}) \quad (1)$$

Bem.: (i) Für $X = Y$ ist $\hat{\sigma}_{X,X} = \hat{\sigma}_X^2$

(ii) $\hat{\sigma}_{X,Y} > 0$ X und Y in der Stichprobe χ positiv korreliert

(iii) $\hat{\sigma}_{X,Y} < 0$ X und Y in der Stichprobe χ negativ korreliert

(iv) $\hat{\sigma}_{X,Y} = \hat{\sigma}_{Y,X}$

Hanno Gottschalk

Stochastik für Informatiker – 9 / 32

Der Korrelationskoeffizient

Die Stärke des Zusammenhanges zwischen X und Y erschließt sich nicht unmittelbar aus der Kovarianz.

Z.B. $x_i \rightarrow ax_i + b \Rightarrow \hat{\sigma}_{X,Y} = a\hat{\sigma}_{X,Y}$

$[\hat{\sigma}_{X,Y}] = [X] \times [Y]$, die Kovarianz ist nicht dimensionslos. . .

Def.: Der emp. Korrelationskoeffizient ist gegeben durch

$$\hat{r}_{X,Y} = \frac{\hat{\sigma}_{X,Y}}{\hat{\sigma}_X \hat{\sigma}_Y} \quad (2)$$

($\hat{\sigma}_X, \hat{\sigma}_Y > 0$ emp. Standardabweichung von X und Y)

Falls $\hat{\sigma}_X = 0$ oder $\hat{\sigma}_Y = 0$ ist $\hat{r}_{X,Y}$ nicht definiert.

Hanno Gottschalk

Stochastik für Informatiker – 10 / 32

Eigenschaften des Korrelationskoeffizienten

Es gilt: (i) $\hat{r}_{X,Y}$ ist dimensionslos

(ii) $-1 \leq \hat{r}_{X,Y} \leq 1$

(iii) $|\hat{r}_{X,Y}| = 1 \Leftrightarrow$ es gibt a, b so dass $y_i = ax_i + b \forall i = 1, \dots, n$

(iv) $x_i \rightarrow ax_i + b$ läßt $\hat{r}_{X,Y}$ unverändert ✓

Denn: (i) $[\hat{\sigma}_X] = [X]$ und $[\hat{\sigma}_Y] = [Y]$

(ii) Nach der Cauchy-Schwarz Ungleichung

$$\left| \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right| \leq \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (3)$$

Also $|(n-1)\hat{\sigma}_{X,Y}| \leq \sqrt{(n-1)}\hat{\sigma}_X \sqrt{(n-1)}\hat{\sigma}_Y$ (iii) In (3) gilt $\Leftrightarrow (y_i - \bar{y}) = a(x_i - \bar{x}) \Rightarrow$ setze $b = \bar{y} - a\bar{x}$

Hanno Gottschalk

Stochastik für Informatiker – 11 / 32

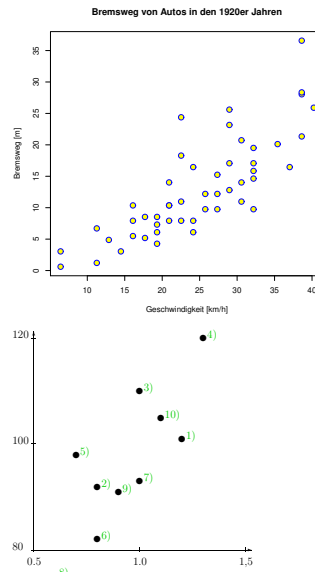
Beispiele für Korrelation

Korrelationskoeffizient Geschwindigkeit
— Bremsweg

$$\hat{r}_{X,Y} = 0.807$$

Korrelationskoeffizient Werbungskosten
— Verkäufe

$$\hat{r}_{X,Y} = 0.875$$

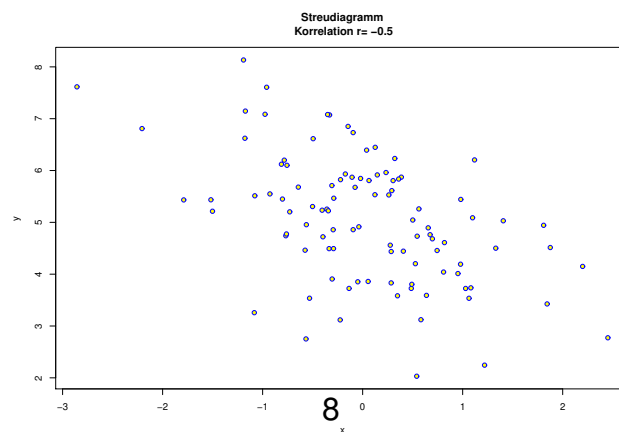
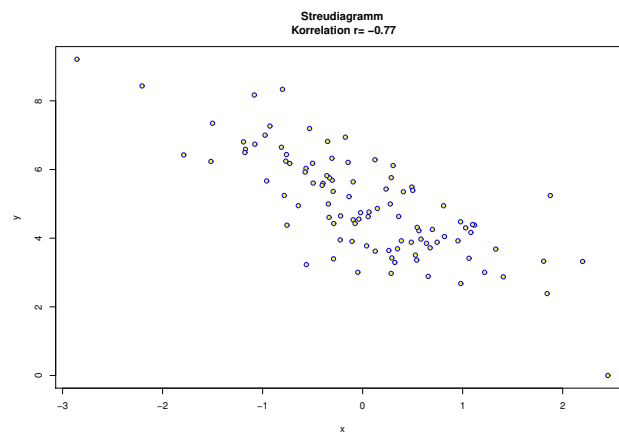
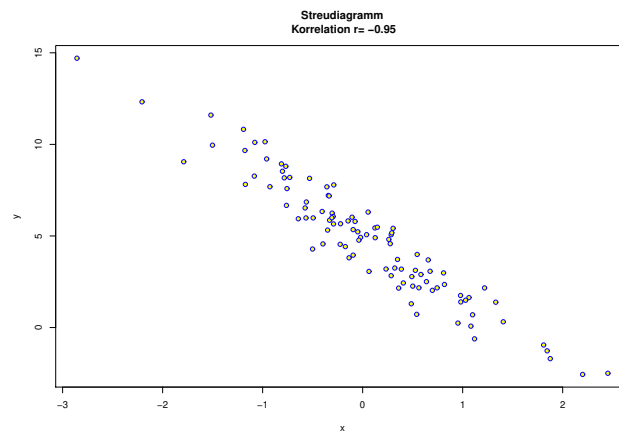
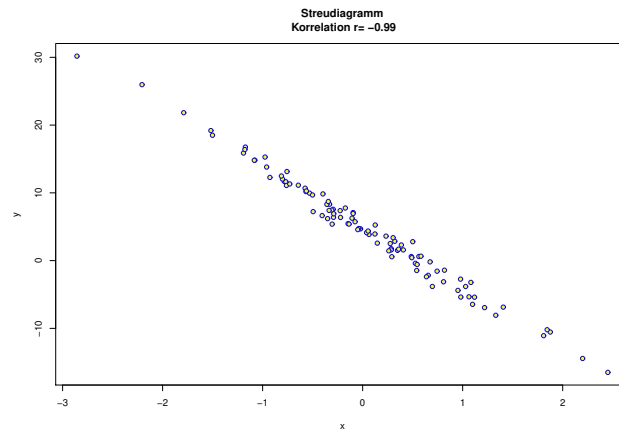


Hanno Gottschalk

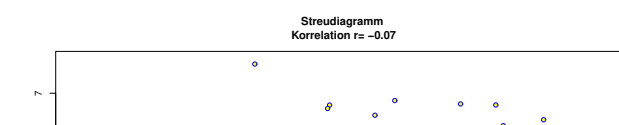
Stochastik für Informatiker – 12 / 32

Einstufung der Korrelation

- $|\hat{r}_{X,Y}| \leq 0.3$ schwache oder keine Korrelation
- $0.3 < |\hat{r}_{X,Y}| \leq 0.9$ mittelstarke Korrelation
- $|\hat{r}_{X,Y}| > 0.9$ starke Korrelation



8



Kausale Abhängigkeiten

Auch der Korrelationskoeffizient ist symmetrisch (wie die Kovarianz)

$$\hat{r}_{X,Y} = \hat{r}_{Y,X}$$

Schon deshalb sagt der Korrelationskoeffizient nichts über kausale Abhängigkeit

Die kausalen Abhängigkeiten ergeben sich aus einer Interpretation, die per se nicht von den Daten her kommt.

- Mehr Werbung führt zu mehr Verkauf ✓
- Mehr Verkauf führt zu mehr Werbung ?

Hanno Gottschalk

Stochastik für Informatiker – 15 / 32

Scheinkorrelation

Es gibt eine verborgene Ursache Z , so dass Zunahme von Z zur Zunahme von X und Zunahme (Abnahme) von Y führt

Dann erscheinen X und Y positiv (negativ) korreliert. . .

Würde man den Z - Effekt herausrechnen, wäre die Korrelation nicht vorhanden oder sogar umgekehrt!

- Verkaufszahlen von Importware und lokaler Produktion sind positiv korreliert
- Je mehr Import, desto mehr lokale Produktion ?
- Die versteckte Ursache Konjunktur treibt beide in dieselbe Richtung

Hanno Gottschalk

Stochastik für Informatiker – 16 / 32

Beispiel Scheinkorrelation

- Medizinische Studie
- Untersuchen Heilerfolg in Abhängigkeit von Dosierung
- Es resultiert eine negative Korrelation $\hat{r}_{\text{Dosis, Erfolg}} = -0.25$
- Medikament Mist ?

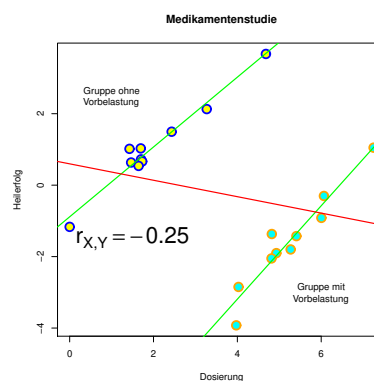
Hanno Gottschalk

Stochastik für Informatiker – 17 / 32



Beispiel Scheinkorrelation II

Blick auf die Daten — immer eine gute Idee:
'let the data speak'!



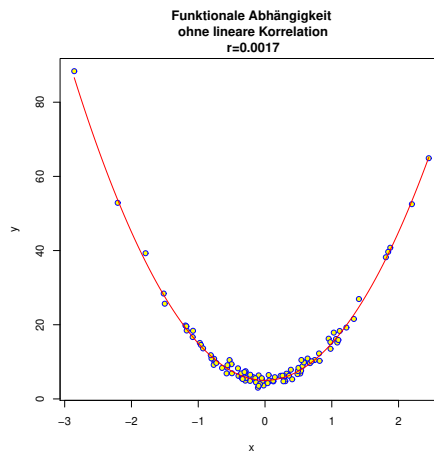
Der Gruppe mit Vorbelastung wurde eine höhere Dosis verabreicht - die Heilaussichten waren für diese Gruppe aber von vornherein geringer!

Hanno Gottschalk

Stochastik für Informatiker – 18 / 32

Abhängigkeit ohne lin. Korrelation

Merkmale X und Y können starke Abhängigkeit aufweisen, ohne linear korreliert zu sein!



Hanno Gottschalk

Stochastik für Informatiker – 19 / 32

Berechnung von Ausgleichsgraden (lin. Reg.)

20 / 32

Vorbemerkung - Modellierung I

Eine Perfekte funktionale Abhängigkeit von Y von X hat die Form

$$Y = f(X) \quad (4)$$

Oft verursachen jedoch noch weitere, nicht gemessene Faktoren $Z, U, V \dots$ eine Veränderung von $Y \dots$

Funktionales Modell der Statistik:

$$Y = f(X) + \text{statistische Schwankungen} \quad (5)$$

$$y_i = f(x_i) + \epsilon_i, \quad i = 1, \dots, n \quad (6)$$

Def.: Die Werte $\epsilon_i = y_i - f(x_i)$ heißen die Residuen (oder Reststreuung) des funktionalen Modells f gegeben die Werte x_1, \dots, x_n und y_1, \dots, y_n

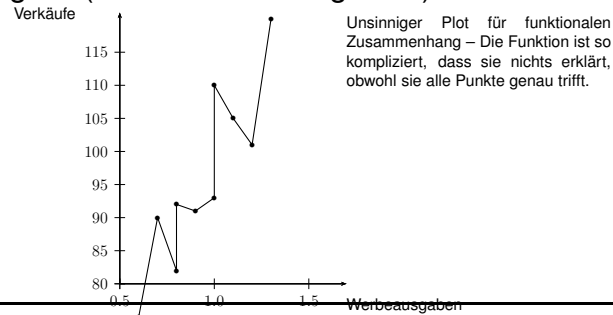
Hanno Gottschalk

Stochastik für Informatiker – 21 / 32

Vorbemerkung - Modellierung II

Aufgabe der Modellbildung:

- Auffinden einer *guten* Funktion f , die Abhängigkeit beschreibt
- einfach zu handhaben
- Erklärungskraft groß (u.a. Reststreuung klein)



Hanno Gottschalk

Stochastik für Informatiker – 22 / 32

Overfitting

Bei verrauschten Daten, vermeide overfitting

overfitting liegt vor, wenn die funktionale Abhängigkeit stark von einzelnen Datenpunkten abhängt

Als Daumenregel: Der Datensatz sollte c.a. $10\times$ größer sein, als die Zahl der freien Parameter des Modells ...

Es sei denn, wir verwenden moderne 'Shrinkage' Methoden zur Modellanpassung (mit Bayes-Motivation)

Genauer kann man nur mit Hilfe von Versuchsplan, Vorwissen über Rauschen in Y etc. sagen

Der Faktor 10 ist kein Dogma!

Hanno Gottschalk

Stochastik für Informatiker – 23 / 32

Lineare Abhängigkeit

Das einfachste Modell für die Abhängigkeit von X und Y ist die *lineare Abhängigkeit*

$$Y = \beta X + \alpha + \text{Reststreuung} \quad (7)$$

Die freien Parameter $\alpha, \beta \in \mathbb{R}$ werden so angepasst, dass die Reststreuung minimal wird.

- einfache Handhabung ✓
- gute Interpretierbarkeit von α und β ✓

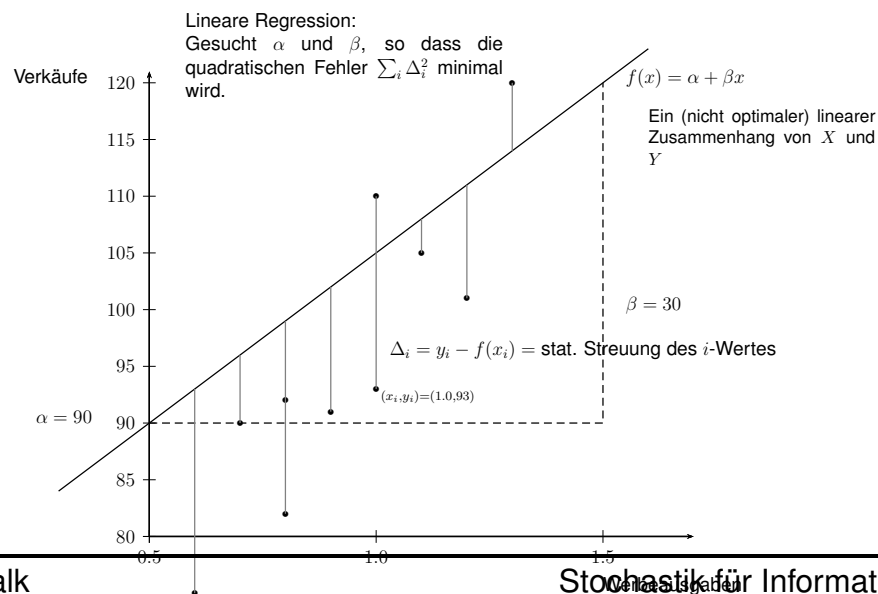
Z.B.: α Verkäufe ohne Werbung, β Ertragssteigerung durch Werbeeinsatz

Wie finde ich die besten Werte $\hat{\alpha}, \hat{\beta}$ für die freien Parameter?

Hanno Gottschalk

Stochastik für Informatiker – 24 / 32

Kleinste Quadrate (least squares)



Hanno Gottschalk

Stochastik für Informatiker – 25 / 32

Lin. Reg. – Berechnung

Minimiere die Fehlerquadrate

$$Q(\alpha, \beta) = \sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 \quad (8)$$

$$0 \stackrel{!}{=} \frac{\partial Q(\alpha, \beta)}{\partial \alpha} = -2 \sum_{i=1}^n (y_i - (\alpha + \beta x_i))$$

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{\partial Q(\alpha, \beta)}{\partial \beta} \\ &= -2 \sum_{i=1}^n (y_i - (\alpha + \beta x_i)) x_i = 2 \sum_{i=1}^n (\alpha x_i + \beta x_i^2 - x_i y_i) \end{aligned}$$

Teile beide Gleichungen durch $n \Rightarrow$

Hanno Gottschalk

Stochastik für Informatiker – 26 / 32

Lin. Reg. – Berechnung

Die Lösung $\hat{\alpha}, \hat{\beta}$ erfüllt also

$$\begin{aligned} 0 &= \hat{\alpha} + \hat{\beta} \bar{x} - \bar{y} \\ 0 &= \hat{\alpha} \bar{x} + \hat{\beta} \overline{x^2} - \overline{xy} \\ &= (\bar{y} - \hat{\beta} \bar{x}) \bar{x} + \hat{\beta} \overline{x^2} - \overline{xy} = \frac{n-1}{n} (\hat{\beta} \hat{\sigma}_X^2 - \hat{\sigma}_{X,Y}) \end{aligned}$$

Satz: Die optimalen Koeffizienten $\hat{\alpha}$ und $\hat{\beta}$ (im Sinne der kleinsten Fehlerquadrate) sind gegeben durch

$$\hat{\beta} = \frac{\hat{\sigma}_{X,Y}}{\hat{\sigma}_X^2} \quad (9)$$

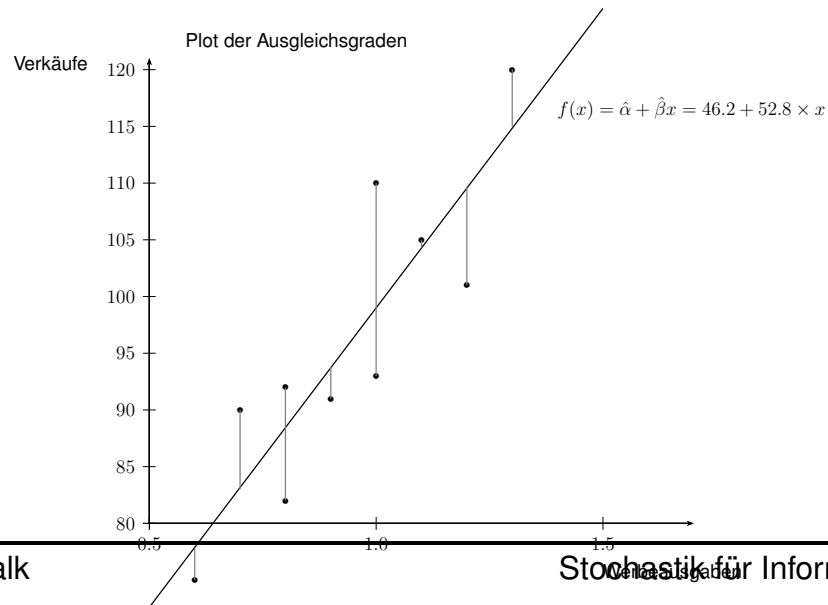
$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad (10)$$

Bemerkung: $[\hat{\beta}] = [Y]/[X]$ und $[\hat{\alpha}] = [Y]$

Hanno Gottschalk

Stochastik für Informatiker – 27 / 32

Ausgleichsgrade Verkäufe



Hanno Gottschalk

Stochastik für Informatiker – 28 / 32

Effektiv Rechnen

29 / 32

Rechenbeispiel

Beispiel:

- Auf Bauernhof werden zwei Futtermittel A und B für Hühner
- A ist billiger und B ist teurer
- X der Anteil des teuren Futtermittels B
- Y Anzahl der gelegten Eier

Merkmal/Wert Nr.	1	2	3	4	5
X =Anteil von Futter B	0	$1/4$	$1/2$	$3/4$	1
Y =gelegte Eier	95	102	104	108	113

Hanno Gottschalk

Stochastik für Informatiker – 30 / 32

Rechnen mit Taschenrechner

	1	2	3	4	5	Σ	$\bar{}$
X	0	1/4	1/2	3/4	1	2.5	0.5
Y	95	102	104	108	113	522	104.4
X^2	0	0.0625	0.25	0.5625	1	1.875	0.375
Y^2	9025	10404	10816	11664	12769	54678	10935.6
XY	0	25.5	52	81	113	271.5	54.3

a) $\bar{x} = 0.5$, $\bar{y} = 104.4$.

$$s_X^2 = \frac{n}{n-1} (\overline{x^2} - \bar{x}^2) = \frac{5}{4} (0.375 - 0.5^2) = 0.15625$$

$$s_Y^2 = \frac{n}{n-1} (\overline{y^2} - \bar{y}^2) = \frac{5}{4} (10935.6 - 104.4^2) = 45.3$$

b) $\text{Cov}_{X,Y} = \frac{n}{n-1} (\overline{xy} - \bar{x}\bar{y}) = \frac{5}{4} (54.3 - 0.5 \times 104.4) = 2.625$

c) $\hat{\beta} = \frac{\text{Cov}_{X,Y}}{s_X^2} = \frac{2.625}{0.15625} = 16.8$, daher $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = 104.4 - 16.8 \times 0.5 = 96 \Rightarrow f(x) = 96 + 16.8 \times x$.

Hanno Gottschalk

Stochastik für Informatiker – 31 / 32

R – wie Rechnen

Mit R ist lineare Regression ein Kinderspiel!

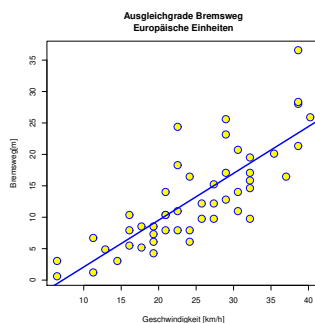
```
R> ?cars # Daten in R vorhanden – US Einheiten!
```

```
R> myModel=lm(data=cars,dist~speed)
# lm für linear model
```

```
R> myModel$coefficients # gefittete Werte
```

```
R> plot(cars) # der Streuplot
```

```
R> abline(myModel) # die Ausgleichsgrade
```



Hanno Gottschalk

Stochastik für Informatiker – 32 / 32