



Allgemeines Lineares Modell

BSc Psychologie SoSe 2022

Prof. Dr. Dirk Ostwald

(1) Regression

Methode der kleinsten Quadrate

Einfache lineare Regression

Selbstkontrollfragen

Methode der kleinsten Quadrate

Einfache lineare Regression

Selbstkontrollfragen

Anwendungsszenario

Psychotherapie



Mehr Therapiestunden

⇒ Höhere Wirksamkeit?

Unabhängige Variable

- Anzahl Therapiestunden

Abhängige Variable

- Symptomreduktion

Methode der kleinsten Quadrate

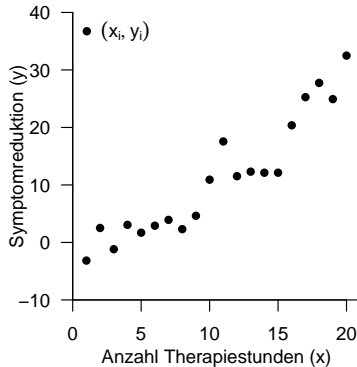
Beispieldatensatz

$i = 1, \dots, 20$ Patient:innen, y_i Symptomreduktion bei Patient:in i , x_i Anzahl Therapiestunden von Patient:in i

y_i	x_i
-3.15	1
2.52	2
-1.18	3
3.06	4
1.70	5
2.91	6
3.92	7
2.31	8
4.63	9
10.91	10
17.56	11
11.52	12
12.31	13
12.12	14
12.13	15
20.37	16
25.26	17
27.75	18
24.93	19
32.49	20

Methode der kleinsten Quadrate

Beispieldatensatz



Welcher funktionaler Zusammenhang zwischen x und y liegt den Daten zugrunde?

Definition (Ausgleichsgerade)

Für $\beta := (\beta_0, \beta_1)^T \in \mathbb{R}^2$ heißt die linear-affine Funktion

$$f_\beta : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto f_\beta(x) := \beta_0 + \beta_1 x, \quad (1)$$

für die für eine Wertemenge $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ die Funktion

$$q : \mathbb{R}^2 \rightarrow \mathbb{R}_{\geq 0}, \beta \mapsto q(\beta) := \sum_{i=1}^n (y_i - f_\beta(x_i))^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad (2)$$

der quadrierten vertikalen Abweichungen der y_i von den Funktionswerten $f_\beta(x_i)$ ihr Minimum annimmt, die *Ausgleichsgerade* für die Wertemenge $\{(x_1, y_1), \dots, (x_n, y_n)\}$.

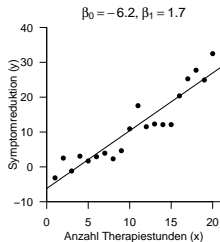
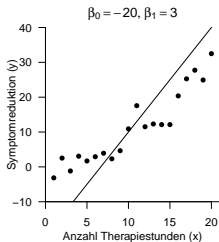
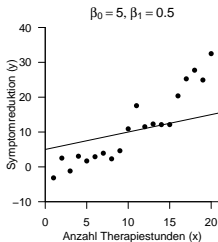
Bemerkungen

- Wir nehmen hier ohne Beweis an, dass das Minimum von q eindeutig ist.

Methode der kleinsten Quadrate

Linear-affine Funktionen $f_{\beta}(x) := \beta_0 + \beta_1 x$

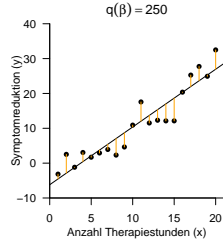
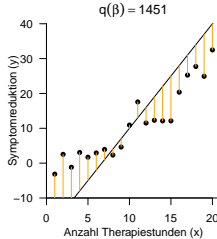
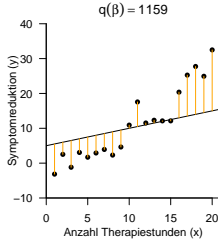
- β_0 : Schnittpunkt von Gerade und y -Achse (“Offset Parameter”)
- β_1 : y -Differenz pro x -Einheitsdifferenz (“Steigungsparameter”)



Methode der kleinsten Quadrate

Funktion der quadrierten vertikalen Abweichungen

$$q(\beta) := \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad (3)$$



— $y_i - (\beta_0 + \beta_1 x_i)$ für $i = 1, \dots, n$

Theorem (Ausgleichsgerade)

Für eine Wertemenge $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ hat die Ausgleichsgerade die Form

$$f_\beta : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto f_\beta(x) := \hat{\beta}_0 + \hat{\beta}_1 x, \quad (4)$$

wobei mit der Stichprobenkovarianz c_{xy} der (x_i, y_i) -Werte, der Stichprobenvarianz s_x^2 der x_i -Werte und den Stichprobenmitteln \bar{x} und \bar{y} der x_i - und y_i -Werte, respektive, gilt, dass

$$\hat{\beta}_1 = \frac{c_{xy}}{s_x^2} \text{ und } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (5)$$

Bemerkungen

- Mit den Definitionen von c_{xy} und s_x^2 gilt also

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

- Man spricht hier von der Stichprobenkovarianz c_{xy} , auch wenn die Werte x_1, \dots, x_n oft nicht als Realisierungen einer Stichprobe X_1, \dots, X_n verstanden werden, sondern als gegebene oder selbst gewählte Zahlen.

Methode der kleinsten Quadrate

Beweis

Wir betrachten die Summe der quadrierten vertikalen Abweichungen der y_i von den Funktionswerten $f(x_i)$ als Funktion von β_0 und β_1 und bestimmen Werte $\hat{\beta}_0$ und $\hat{\beta}_1$, für die diese Funktion ihr Minimum annimmt, die Summe der quadrierten vertikalen Abweichungen der y_i von den Funktionswerten $f(x_i)$ also minimal ist. Wir betrachten also die Funktion

$$q : \mathbb{R}^2 \rightarrow \mathbb{R}, (\beta_0, \beta_1) \mapsto q(\beta_0, \beta_1) := \sum_{i=1}^n \left(y_i - (\beta_0 + \beta_1 x_i) \right)^2. \quad (7)$$

Um das Minimum dieser Funktion zu bestimmen, berechnen wir zunächst die partiellen Ableitungen hinsichtlich β_0 und β_1 und setzen diese gleich 0. Es ergibt sich zunächst

$$\begin{aligned} \frac{\partial}{\partial \beta_0} q(\beta_0, \beta_1) &= \frac{\partial}{\partial \beta_0} \left(\sum_{i=1}^n \left(y_i - (\beta_0 + \beta_1 x_i) \right)^2 \right) \\ &= \sum_{i=1}^n \frac{\partial}{\partial \beta_0} \left(y_i - (\beta_0 + \beta_1 x_i) \right)^2 \\ &= \sum_{i=1}^n 2 \left(y_i - (\beta_0 + \beta_1 x_i) \right) \frac{\partial}{\partial \beta_0} \left(y_i - \beta_0 - \beta_1 x_i \right) \\ &= -2 \sum_{i=1}^n \left(y_i - \beta_0 - \beta_1 x_i \right) \end{aligned} \quad (8)$$

Methode der kleinsten Quadrate

Beweis (fortgeführt)

Weiterhin ergibt sich

$$\begin{aligned}\frac{\partial}{\partial \beta_1} q(\beta_0, \beta_1) &= \frac{\partial}{\partial \beta_1} \left(\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \right) \\&= \sum_{i=1}^n \frac{\partial}{\partial \beta_1} (y_i - (\beta_0 + \beta_1 x_i))^2 \\&= \sum_{i=1}^n 2 (y_i - (\beta_0 + \beta_1 x_i)) \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i) \\&= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i\end{aligned}\tag{9}$$

Nullsetzen beider partieller Ableitungen ergibt dann

$$\begin{aligned}\frac{\partial}{\partial \beta_0} q(\beta_0, \beta_1) &= 0 \text{ und } \frac{\partial}{\partial \beta_1} q(\beta_0, \beta_1) = 0 \\ \Leftrightarrow -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) &= 0 \text{ und } -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0 \\ \Leftrightarrow \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) &= 0 \text{ und } \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0\end{aligned}\tag{10}$$

Beweis (fortgeführt)

und weiter

$$\begin{aligned} \sum_{i=1}^n y_i - \sum_{i=1}^n \beta_0 - \beta_1 \sum_{i=1}^n x_i = 0 \text{ und } \sum_{i=1}^n y_i x_i - \sum_{i=1}^n \beta_0 x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0 \\ \Leftrightarrow \beta_0 n + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \text{ und } \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i \end{aligned} \quad (11)$$

Das sich hier ergebende Gleichungssystem

$$\begin{aligned} \beta_0 n + \beta_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n y_i x_i \end{aligned} \quad (12)$$

wird *System der Normalgleichungen* genannt und beschreibt die notwendige Bedingung für ein Minimum von q . Auflösen dieses Gleichungssystems nach β_0 und β_1 liefert dann die Werte $\hat{\beta}_0$ und $\hat{\beta}_1$ des Theorems.

Methode der kleinsten Quadrate

Beweis (fortgeführt)

Um dies zu sehen, halten wir zunächst fest, dass mit der ersten Gleichung des Systems der Normalgleichungen gilt

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \Leftrightarrow \hat{\beta}_0 + \hat{\beta}_1 \bar{x} = \bar{y} \Leftrightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (13)$$

Einsetzen der Form von $\hat{\beta}_0$ in die zweite Gleichung des Systems der Normalgleichungen ergibt dann zunächst

$$\begin{aligned} & \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i \\ & \Leftrightarrow (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i \\ & \Leftrightarrow \bar{y} \sum_{i=1}^n x_i - \hat{\beta}_1 \bar{x} \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i \quad (14) \\ & \Leftrightarrow -\hat{\beta}_1 \bar{x} \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i \\ & \Leftrightarrow \hat{\beta}_1 \left(\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i \right) = \sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i \end{aligned}$$

Beweis (fortgeführt)

Wir halten nun zunächst fest, dass gilt

$$\begin{aligned}\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i &= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \bar{x} \sum_{i=1}^n x_i \\&= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n \left(\frac{1}{n} \sum_{i=1}^n x_i \right) \bar{x} \\&= \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2 \\&= \sum_{i=1}^n \left(x_i^2 - 2\bar{x}x_i + \bar{x}^2 \right) \\&= \sum_{i=1}^n \left(x_i - \bar{x} \right)^2.\end{aligned}\tag{15}$$

Beweis (fortgeführt)

Weiterhin halten wir zunächst fest, dass gilt

$$\begin{aligned}\sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i - n \bar{y} \bar{x} + n \bar{y} \bar{x} \\&= \sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i - \sum_{i=1}^n y_i \bar{x} + \sum_{i=1}^n \bar{y} \bar{x} \\&= \sum_{i=1}^n y_i x_i - \sum_{i=1}^n y_i \bar{x} - \sum_{i=1}^n \bar{y} x_i + \sum_{i=1}^n \bar{y} \bar{x} \\&= \sum_{i=1}^n \left(y_i x_i - y_i \bar{x} - \bar{y} x_i + \bar{y} \bar{x} \right) \\&= \sum_{i=1}^n \left(y_i - \bar{y} \right) \left(x_i - \bar{x} \right).\end{aligned}\tag{16}$$

Beweis (fortgeführt)

In der Fortsetzung von (14) ergibt sich dann

$$\begin{aligned}\hat{\beta}_1 \left(\sum_{i=1}^n x_i^2 - \bar{x} \sum_{i=1}^n x_i \right) &= \sum_{i=1}^n y_i x_i - \bar{y} \sum_{i=1}^n x_i \\ \Leftrightarrow \hat{\beta}_1 \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) &= \sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x}) \\ \Leftrightarrow \hat{\beta}_1 &= \frac{\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ \Leftrightarrow \hat{\beta}_1 &= \frac{c_{xy}}{s_x^2}.\end{aligned}\tag{17}$$

□

Beispieldatensatz Analyse

```
# Einlesen des Beispieldatensatzes
fname      = file.path(getwd(), "1_Daten", "1_Regression.csv")
D          = read.table(fname, sep = ",", header = TRUE)

# Stichprobenstatistiken
x_bar      = mean(D$x_i)           # Stichprobenmittel der x_i-Werte
y_bar      = mean(D$y_i)           # Stichprobenmittel der y_i-Werte
s2x        = var(D$x_i)            # Stichprobenvarianz der x_i-Werte
cxy        = cov(D$x_i, D$y_i)     # Stichprobenkovarianz der (x_i, y_i)-Werte

# Ausgleichsgeradenparameter
beta_1_hat = cxy/s2x               # \hat{\beta}_1, Steigungsparameter
beta_0_hat = y_bar - beta_1_hat*x_bar # \hat{\beta}_0, Offset Parameter

# Ausgabe
cat("beta_0_hat:", beta_0_hat,
    "\nbeta_1_hat:", beta_1_hat)

> beta_0_hat: -6.19
> beta_1_hat: 1.66
```

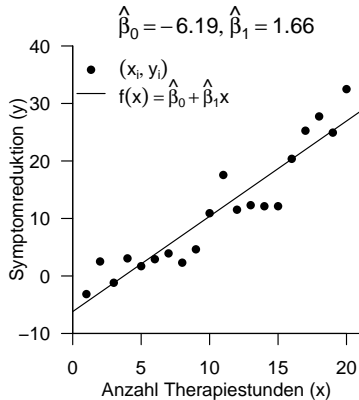
Beispieldatensatz Visualisierung

```
# Datenwerte
plot(
  D$x_i,
  D$y_i,
  pch      = 16,
  xlab     = "Anzahl Therapiestunden (x)",
  ylab     = "Symptomreduktion (y)",
  xlim     = c(0,21),
  ylim     = c(-10, 40),
  main     = TeX("\\hat{\\beta}_0 = -6.19, \\hat{\\beta}_1 = 1.66$"))

# Ausgleichsgerade
abline(
  coef     = c(beta_0_hat, beta_1_hat),
  lty      = 1,
  col      = "black")

# Legende
legend(
  "topleft",
  c(TeX("\\$(x_i,y_i)$"), TeX("$f(x) = \\hat{\\beta}_0 + \\hat{\\beta}_1x$")),
  lty      = c(0,1),
  pch      = c(16, NA),
  bty      = "n")
```

Beispieldatensatz Visualisierung



Definition (Ausgleichspolynom)

Für $\beta := (\beta_0, \dots, \beta_k)^T \in \mathbb{R}^{k+1}$ heißt die Polynomfunktion k ten Grades

$$f_\beta : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto f_\beta(x) := \sum_{i=0}^k \beta_i x^i, \quad (18)$$

für die für eine Wertemenge $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^2$ die Funktion

$$q : \mathbb{R}^2 \rightarrow \mathbb{R}_{\geq 0}, \beta \mapsto q(\beta) := \sum_{i=1}^n \left(y_i - f_\beta(x_i) \right)^2 = \sum_{i=1}^n \left(y_i - \sum_{i=0}^k \beta_i x^i \right)^2 \quad (19)$$

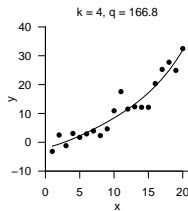
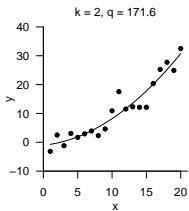
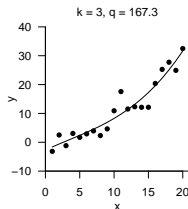
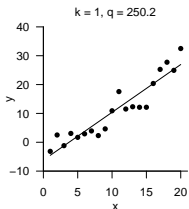
der quadrierten vertikalen Abweichungen der y_i von den Funktionswerten $f_\beta(x_i)$ ihr Minimum annimmt, das *Ausgleichspolynom* k ten Grades für die Wertemenge $\{(x_1, y_1), \dots, (x_n, y_n)\}$.

Bemerkungen

- Wir nehmen hier ohne Beweis an, dass das Minimum von q eindeutig ist.
- Die Ausgleichsgerade ist das Ausgleichspolynom ersten Grades.
- Die Parameterwerte $\hat{\beta}_0, \dots, \hat{\beta}_k$ für die q bei gegebener Wertemenge ihr Minimum annehmen werden an späterer Stelle im Rahmen der Theorie des Allgemeinen Linearen Modells ALMs bestimmt werden.

Methode der kleinsten Quadrate

Beispieldatensatz Ausgleichspolynome 1ten bis 4ten Grades



$$\bullet (x_i, y_i) \quad \text{---} \quad f_{\hat{\beta}}(x) = \sum_{i=0}^k \hat{\beta}_i x^i$$

Methode der kleinsten Quadrate

Einfache lineare Regression

Selbstkontrollfragen

Motivation

Eine Ausgleichsgerade erlaubt Aussagen über unbeobachtete y Werte für x Werte. Der Wert von $q(\hat{\beta})$ quantifiziert die Güte der Ausgleichsgeradenpassung. Eine Ausgleichsgerade erlaubt allerdings nur implizite Aussagen über die mit der Anpassung verbundene Unsicherheit.

In der einfachen linearen Regression wird die Idee einer Ausgleichsgerade um eine probabilistische Komponente (normalverteilte Fehlervariable) erweitert, um quantitative Aussagen über die mit einer Ausgleichsgeradenanpassung verbundene Unsicherheit machen zu können. Weiterhin erlaubt die einfache lineare Regression, einen Hypothesentest- basierten Zugang zur Einschätzung der angepassten Parameterwerte $\hat{\beta}_0$ und $\hat{\beta}_1$.

Wir betrachten hier zunächst nur das probabilistische Modell der einfachen linearen Regression sowie die auf ihm basierende Maximum Likelihood Schätzung der Parameter β_0 und β_1 . Die Bewertung von Parameterschätzerunsicherheit sowie parameterzentrierte Hypothesentests behandeln wir an späterer Stelle zunächst im Allgemeinen.

Definition (Generatives Modell der einfachen linearen Regression)

Für $i = 1, \dots, n$ sei

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (20)$$

wobei

- $x_i \in \mathbb{R}$ fest vorgegebene sogenannte *Prädiktorwerte* oder *Regressorwerte* sind,
- $\beta_0, \beta_1 \in \mathbb{R}$ wahre, aber unbekannte, Parameterwerte sind und
- $\varepsilon_i \sim N(0, \sigma^2)$ unabhängige und identisch normalverteilte nicht-beobachtbare Zufallsvariablen mit wahrem, aber unbekanntem, Parameter $\sigma^2 > 0$ sind.

Dann heißt (20) *Generatives Modell der einfachen linearen Regression*.

Bemerkungen

- Das Modell der einfachen linearen Regression hat drei Parameter, $\beta_0, \beta_1 \in \mathbb{R}$ und $\sigma^2 > 0$.

Theorem (Normalverteilungsmodell der einfachen linearen Regression)

Das generative Modell der einfachen linearen Regression

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \text{ mit } \varepsilon_i \sim N(0, \sigma^2) \text{ u.i.v. für } i = 1, \dots, n \quad (21)$$

lässt sich äquivalent in der Form

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \text{ u.i.v. für } i = 1, \dots, n \quad (22)$$

schreiben.

Bemerkungen

- Wir bezeichnen

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \text{ u.i.v. für } i = 1, \dots, n \quad (23)$$

als *Normalverteilungsmodell der einfachen linearen Regression*.

Einfache lineare Regression

Beweis

Wir zeigen die Äquivalenz für ein i , die Unabhängigkeit der Y_i zeigen wir an späterer Stelle im Rahmen des Allgemeinen Linearen Modells. Die Äquivalenz beider Modellformen für ein i folgt direkt aus der Transformation normalverteilter Zufallsvariablen durch linear-affine Funktionen (cf. (8) Transformationen der Normalverteilung). Speziell gilt im vorliegenden Fall für $\varepsilon_i \sim N(0, \sigma^2)$, dass

$$Y_i = f(\varepsilon_i) \text{ mit } f: \mathbb{R} \rightarrow \mathbb{R}, \varepsilon_i \mapsto f(\varepsilon_i) := \varepsilon_i + (\beta_0 + \beta_1 x_i) \quad (24)$$

Mit dem WDF Transformationstheorem bei linear-affinen Abbildungen folgt dann

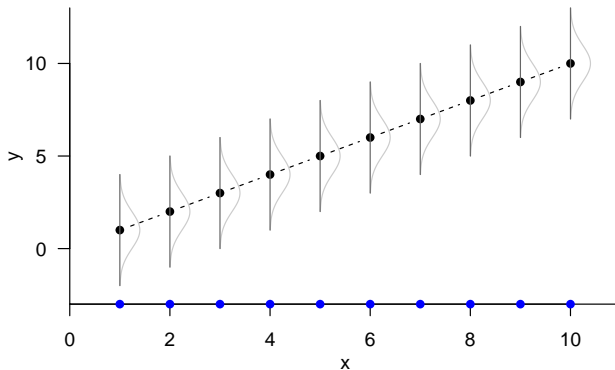
$$\begin{aligned} p_{Y_i}(y_i) &= \frac{1}{|1|} p_{\varepsilon_i} \left(\frac{y_i - \beta_0 - \beta_1 x_i}{1} \right) \\ &= N \left(x_i - \beta_0 - \beta_1 x_i; 0, \sigma^2 \right) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2\sigma^2} (x_i - \beta_0 - \beta_1 x_i - 0)^2 \right) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2\sigma^2} (x_i - (\beta_0 + \beta_1 x_i))^2 \right) \\ &= N \left(x_i; \beta_0 + \beta_1 x_i, \sigma^2 \right), \end{aligned} \quad (25)$$

also

$$Y_i \sim N \left(\beta_0 + \beta_1 x_i, \sigma^2 \right). \quad (26)$$

Einfache lineare Regression

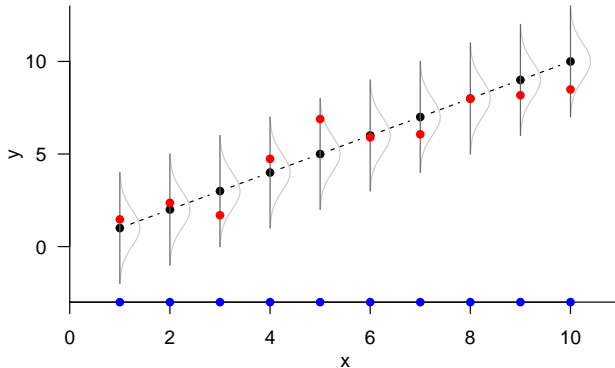
Modell der einfachen linearen Regression



• x_i • $\beta_0 + \beta_1 x_i$ für $\beta_0 := 0, \beta_1 := 1$ — $N(y_i; \beta_0 + \beta_1 x_i, \sigma^2)$ für $\sigma^2 := 1$.

Einfache lineare Regression

Realisierung des Modells der einfachen linearen Regression



• x_i • $\beta_0 + \beta_1 x_i$ für $\beta_0 := 0, \beta_1 := 1$ — $N(y_i; \beta_0 + \beta_1 x_i, \sigma^2)$ für $\sigma^2 := 1$ • (x_i, y_i)

Theorem (Maximum Likelihood Schätzung)

Es sei

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \text{ mit } \varepsilon_i \sim N(0, \sigma^2) \text{ u.i.v. für } i = 1, \dots, n \quad (27)$$

das Modell der einfachen linearen Regression. Dann sind Maximum Likelihood Schätzer der Modellparameter β_0 , β_1 und σ^2 gegeben durch

$$\hat{\beta}_1 := \frac{c_{xy}}{s_x^2}, \quad \hat{\beta}_0 := \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{und} \quad \hat{\sigma}^2 := \frac{1}{n} \sum_{i=1}^n \left(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2. \quad (28)$$

Bemerkungen

- Wir verzichten hier aus Gründen der Übersichtlichkeit auf die ^{ML} Superskripte.
- Die ML Schätzer für β_0 und β_1 sind offenbar mit den Ausgleichsgeradenparametern identisch.

Einfache lineare Regression

Beweis

Wir zeigen zunächst, dass die Ausgleichsgeradenparameter $\hat{\beta}_0$ und $\hat{\beta}_1$ den entsprechenden ML Schätzern gleichen. Dazu halten wir zunächst fest, dass aufgrund der Unabhängigkeit der Y_1, \dots, Y_n die Likelihood-Funktion des Modells der einfachen linearen Regression bezüglich β_0 und β_1 die Form

$$\begin{aligned} L : \mathbb{R}^2 \rightarrow \mathbb{R}_{>0}, (\beta_0, \beta_1) &\mapsto L(\beta_0, \beta_1) := \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - (\beta_0 + \beta_1 x_i))^2\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2\right) \end{aligned} \quad (29)$$

Weil für die Exponentialfunktion gilt, dass für $a < b \leq 0$ gilt, dass $\exp(a) < \exp(b)$ wird der Exponentialterm dieser Likelihood-Funktion maximal, wenn der Term

$$q := \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \geq 0 \quad (30)$$

minimal und damit $-q$ maximal wird. Im Rahmen des Beweises der Ausgleichsgeradenform haben wir aber schon gezeigt, dass der Term (30) für

$$\hat{\beta}_1 := \frac{c_{xy}}{s_x^2} \text{ und } \hat{\beta}_0 := \bar{y} - \hat{\beta}_1 \bar{x} \quad (31)$$

minimal wird, und damit $\hat{\beta}_1$ und $\hat{\beta}_0$ die Likelihood-Funktion maximieren.

Einfache lineare Regression

Beweis (fortgeführt)

In einem zweiten Schritt betrachten wir nun die Likelihood-Funktion des Modells der einfachen linearen Regression bezüglich σ^2 an der Stelle von $\hat{\beta}_0$ und $\hat{\beta}_1$. Wir erhalten die Likelihood-Funktion

$$L : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}, \sigma^2 \mapsto L(\sigma^2) = \left(2\pi\sigma^2\right)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2\right) \quad (32)$$

und die entsprechende Log-Likelihood-Funktion

$$\ell : \mathbb{R}_{>0} \rightarrow \mathbb{R}, \sigma^2 \mapsto \ell(\sigma^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 \quad (33)$$

In Analogie zu der Herleitung des ML Schätzers für σ^2 im Normalverteilungsmodell (cf. (10) Parameterschätzung) ergibt sich unter Beachtung von

$$\hat{\mu}_n^{\text{ML}} = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (34)$$

dann hier

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2. \quad (35)$$

Beispieldatensatz Parameterschätzung

```
# Einlesen des Beispieldatensatzes
fname      = file.path(getwd(), "1_Daten", "1_Regression.csv")
D          = read.table(fname, sep = ",", header = TRUE)

# Stichprobenstatistiken
n          = length(D$y_i)                # Anzahl Datenpunkte
x_bar      = mean(D$x_i)                  # Stichprobenmittel der x_i-Werte
y_bar      = mean(D$y_i)                  # Stichprobenmittel der y_i-Werte
s2x        = var(D$x_i)                   # Stichprobenvarianz der x_i-Werte
cxy        = cov(D$x_i, D$y_i)            # Stichprobenkovarianz der (x_i, y_i)-Werte

# Parameterschätzer
beta_1_hat = cxy/s2x                      # \hat{\beta}_1, Steigungsparameter
beta_0_hat = y_bar - beta_1_hat*x_bar      # \hat{\beta}_0, Offset Parameter
sigsqr_hat = (1/n)*sum((D$y_i-(beta_0_hat+beta_1_hat*D$x_i))^2) # Varianzparameter

# Ausgabe
cat("beta_0_hat:" , beta_0_hat,
    "\nbeta_1_hat:", beta_1_hat,
    "\nsigsqr_hat:", sqrt(sigsqr_hat))

> beta_0_hat: -6.19
> beta_1_hat: 1.66
> sigsq_hat: 3.54
```

Beispieldatensatz Analyse mit `lm()`

```
# Einlesen des Beispieldatensatzes
```

```
library(car)
```

```
> Lade nötiges Paket: carData
```

```
fname      = file.path(getwd(), "1_Daten", "1_Regression.csv")
```

```
D          = read.table(fname, sep = ",", header = TRUE)
```

```
# Analyse mit lm()
```

```
model      = lm(formula = D$y_i ~ D$x_i, data = D)
```

```
print(model)
```

```
>
```

```
> Call:
```

```
> lm(formula = D$y_i ~ D$x_i, data = D)
```

```
>
```

```
> Coefficients:
```

```
> (Intercept)      D$x_i
```

```
>      -6.19        1.66
```

Methode der kleinsten Quadrate

Einfache lineare Regression

Selbstkontrollfragen

References
