



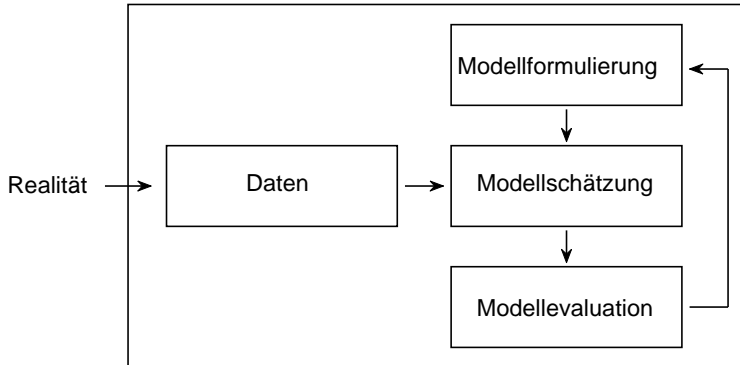
Allgemeines Lineares Modell

BSc Psychologie, SoSe 2024

Joram Soch

(5) Modellformulierung

Naturwissenschaft



Modellformulierung

$$y = X\beta + \varepsilon, \quad \varepsilon \sim N(0_n, \sigma^2 I_n) \quad (1)$$

Modellschätzung

$$\hat{\beta} = (X^T X)^{-1} X^T y, \quad \hat{\sigma}^2 = \frac{1}{n-p} (y - X\hat{\beta})^T (y - X\hat{\beta}) \quad (2)$$

Modellevaluation

$$T = \frac{c^T \hat{\beta} - c^T \beta_0}{\sqrt{\hat{\sigma}^2 c^T (X^T X)^{-1} c}}, \quad F = \frac{(\hat{\varepsilon}_0^T \hat{\varepsilon}_0 - \hat{\varepsilon}^T \hat{\varepsilon})/p_1}{\hat{\varepsilon}^T \hat{\varepsilon}/(n-p)} \quad (3)$$

Standardprobleme Frequentistischer Inferenz

(1) Parameterschätzung

Ziel der Parameterschätzung ist es, einen möglichst guten Tipp für wahre, aber unbekannte, Parameterwerte oder Funktionen dieser abzugeben, typischerweise mithilfe von Daten.

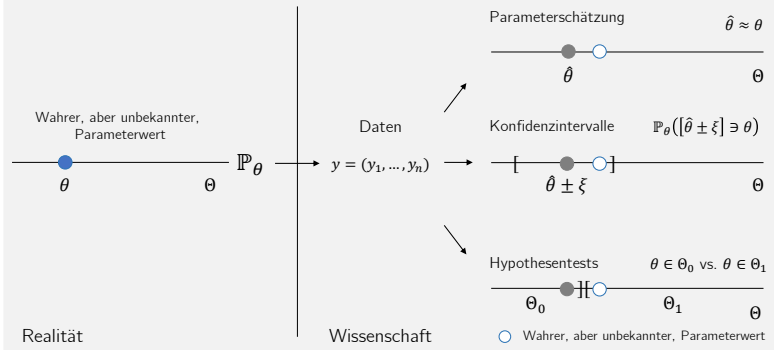
(2) Konfidenzintervalle

Ziel der Bestimmung von Konfidenzintervallen ist es, basierend auf der angenommenen Verteilung der Daten eine quantitative Aussage über die mit Schätzwerten assoziierte Unsicherheit zu treffen.

(3) Hypothesentests

Ziel des Hypothesentestens ist es, basierend auf der angenommenen Verteilung der Daten in einer möglichst zuverlässigen Form zu entscheiden, ob ein wahrer, aber unbekannter Parameterwert in einer von zwei sich gegenseitig ausschließenden Untermengen des Parameterraumes liegt.

Modell und Standardprobleme Frequentistischer Inferenz



$$\theta := (\beta, \sigma^2), \quad \Theta := \mathbb{R}^p \times \mathbb{R}_{>0}, \quad \mathbb{P}_\theta(y) := \mathbb{P}_{\beta, \sigma^2}(y) \quad \text{mit WDF} \quad p_{\beta, \sigma^2}(y) := N(y; X\beta, \sigma^2 I_n)$$

Standardannahmen Frequentistischer Inferenz

Gegeben sei das Allgemeine Lineare Modell. Es wird angenommen, dass ein vorliegender Datensatz eine der möglichen Realisierungen der Daten des Modells ist. Aus Frequentistischer Sicht kann man unendlich oft Datensätze aus einem Modell generieren und zu jedem Datensatz Schätzer oder Statistiken auswerten, z.B. den Betaparameterschätzer:

$$\text{Datensatz (1)} : y^{(1)} = \left(y_1^{(1)}, y_2^{(1)}, \dots, y_n^{(1)} \right)^T \quad \text{mit} \quad \hat{\beta}^{(1)} = (X^T X)^{-1} X^T y^{(1)}$$

$$\text{Datensatz (2)} : y^{(2)} = \left(y_1^{(2)}, y_2^{(2)}, \dots, y_n^{(2)} \right)^T \quad \text{mit} \quad \hat{\beta}^{(2)} = (X^T X)^{-1} X^T y^{(2)}$$

$$\text{Datensatz (3)} : y^{(3)} = \left(y_1^{(3)}, y_2^{(3)}, \dots, y_n^{(3)} \right)^T \quad \text{mit} \quad \hat{\beta}^{(3)} = (X^T X)^{-1} X^T y^{(3)}$$

$$\text{Datensatz (4)} : y^{(4)} = \left(y_1^{(4)}, y_2^{(4)}, \dots, y_n^{(4)} \right)^T \quad \text{mit} \quad \hat{\beta}^{(4)} = (X^T X)^{-1} X^T y^{(4)}$$

$$\text{Datensatz (5)} : y^{(5)} = \dots$$

Um die Qualität statistischer Methoden zu beurteilen betrachtet die Frequentistische Statistik die Wahrscheinlichkeitsverteilungen von Schätzern und Statistiken unter Annahme der Datenverteilung. Was zum Beispiel ist die Verteilung von $\hat{\beta}^{(1)}, \hat{\beta}^{(2)}, \hat{\beta}^{(3)}, \hat{\beta}^{(4)}, \dots$, also die Verteilung der Zufallsvariable $\hat{\beta} := (X^T X)^{-1} X^T y$? Wenn eine statistische Methode im Sinne der Frequentistischen Standardannahmen "gut" ist, dann heißt das also, dass sie bei häufiger Anwendung "im Mittel gut" ist. Im Einzelfall, also im Normalfall nur eines vorliegenden Datensatzes, kann sie auch "schlecht" sein.

Anwendungsbeispiele in Einheiten (5) – (8)

- Einstichproben-T-Test
- Einfache lineare Regression

Anwendungsbeispiele in Einheiten (9) – (14)

- Zweistichproben-T-Tests
- Einfaktorielle Varianzanalyse
- Zweifaktorielle Varianzanalyse
- Multiple Regression
- Kovarianzanalyse

Allgemeine Theorie

Unabhängige und identisch normalverteilte Zufallsvariablen

Einfache lineare Regression

Selbstkontrollfragen

Allgemeine Theorie

Unabhängige und identisch normalverteilte Zufallsvariablen

Einfache lineare Regression

Selbstkontrollfragen

Definition (Allgemeines Lineares Modell)

Es sei

$$y = X\beta + \varepsilon, \quad (4)$$

wobei

- y ein n -dimensionaler beobachtbarer Zufallsvektor ist, der *Daten* genannt wird,
- $X \in \mathbb{R}^{n \times p}$ mit $n > p$ eine vorgegebene Matrix ist, die *Designmatrix* genannt wird,
- $\beta \in \mathbb{R}^p$ ein unbekannter Parametervektor ist, der *Betaparametervektor* genannt wird,
- ε ein n -dimensionaler nicht-beobachtbarer Zufallsvektor ist, der *Zufallsfehler* genannt wird und für den angenommen wird, dass mit einem unbekanntem Varianzparameter $\sigma^2 > 0$ gilt, dass

$$\varepsilon \sim N(0_n, \sigma^2 I_n). \quad (5)$$

Dann heißt (4) *Allgemeines Lineares Modell (ALM)*.

Bemerkungen

- Wir bezeichnen hier mit y sowohl den Zufallsvektor, über dessen Wahrscheinlichkeitsverteilung wir Aussagen machen können (vormals v), als auch eine konkrete Realisierung dieses Zufallsvektors, die einen gegebenen Datensatz darstellt. Es wird also ab hier nicht mehr zwischen dem Zufallsvektor v und seiner Realisierung y unterschieden. In der Literatur wird das ALM in der überwiegenden Zahl der Fälle als $y = X\beta + \varepsilon$ beschrieben.

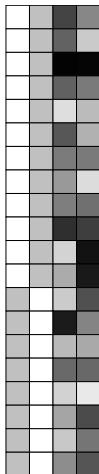
Bemerkungen (fortgeführt)

- Wir nehmen durchgängig an, dass $X \in \mathbb{R}^{n \times p}$ vollen Spaltenrang hat, also dass $\text{rg}(X) = p$.
- y ist ein Zufallsvektor, weil er aus der Addition des Zufallsvektors $\varepsilon \in \mathbb{R}^n$ zu dem Vektor $X\beta \in \mathbb{R}^n$ resultiert.
- Wir nennen $X\beta$ den *deterministischen Modellaspekt* und ε den *probabilistischen Modellaspekt*.
- $n \in \mathbb{N}$ bezeichnet durchgängig die Anzahl an Datenpunkten.
- $p \in \mathbb{N}$ bezeichnet durchgängig die Anzahl an Betaparametern.
- Die Gesamtzahl an Parametern des ALMs ist $p + 1$ (p Betaparameterkomponenten und 1 Varianzparameter).
- Der Betaparametervektor wird auch *Gewichtsvektor*, *Effektvektor* oder *Regressionskoeffizienten* genannt.
- Weil der Kovarianzmatrixparameter von ε als sphärisch angenommen wird, sind die $\varepsilon_1, \dots, \varepsilon_n$ unabhängige normalverteilte Zufallsvariablen mit identischem Varianzparameter. Weil zusätzlich der Erwartungswertparameter von ε als 0_n angenommen wird, sind die $\varepsilon_1, \dots, \varepsilon_n$ auch identisch normalverteilte Zufallsvariablen.
- Für jede Komponente $y_i, i = 1, \dots, n$ von y impliziert (4) nach Definition des Matrixprodukts, dass

$$y_i = \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i = x_{i1}\beta_1 + \dots + x_{ip}\beta_p + \varepsilon_i \quad \text{mit} \quad \varepsilon_i \sim N(0, \sigma^2), \quad (6)$$

wobei $x_{ij} \in \mathbb{R}$ das (ij) te Element der Designmatrix X bezeichnet.

Darstellung einer Designmatrix als Graustufen-Matrixplot (hier: $n = 20$, $p = 4$)



Darstellung einer Designmatrix als Graustufen-Matrixplot (hier: $n = 20$, $p = 4$)

$$X := \begin{pmatrix} x_1 & x_2 & x_3 & x_4 \end{pmatrix} \in \mathbb{R}^{n \times 4} \quad (7)$$

$$\beta := \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} \in \mathbb{R}^4 \quad (8)$$

$$\begin{aligned} X\beta &= \begin{pmatrix} x_1 & x_2 & x_3 & x_4 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} \\ &= \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \in \mathbb{R}^{n \times 1} \end{aligned} \quad (9)$$

Theorem (Datenverteilung des Allgemeinen Linearen Modells)

Es sei

$$y = X\beta + \varepsilon \quad \text{mit} \quad \varepsilon \sim N(0_n, \sigma^2 I_n) \quad (10)$$

das ALM. Dann gilt

$$y \sim N(\mu, \sigma^2 I_n) \quad \text{mit} \quad \mu := X\beta \in \mathbb{R}^n. \quad (11)$$

Beweis

Das Theorem zur linear-affinen Transformation multivariater Normalverteilungen besagt:

$$\xi \in \mathbb{R}^n, \quad \xi \sim N(\mu, \Sigma), \quad A \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m \Rightarrow v := A\xi + b \sim N(A\mu + b, A\Sigma A^T). \quad (12)$$

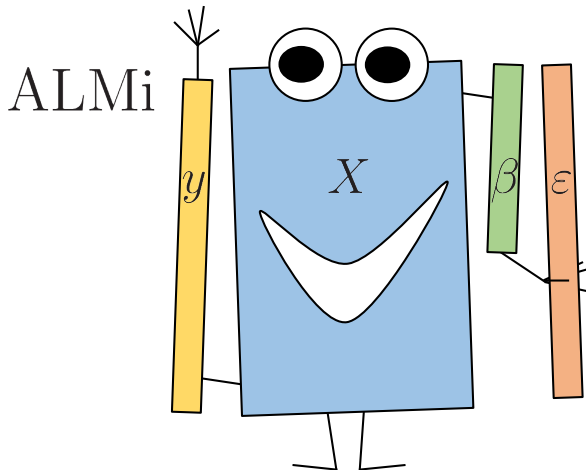
Damit gilt für $\varepsilon \sim N(0_n, \sigma^2 I_n)$ und $y := I_n \varepsilon + X\beta$, dass

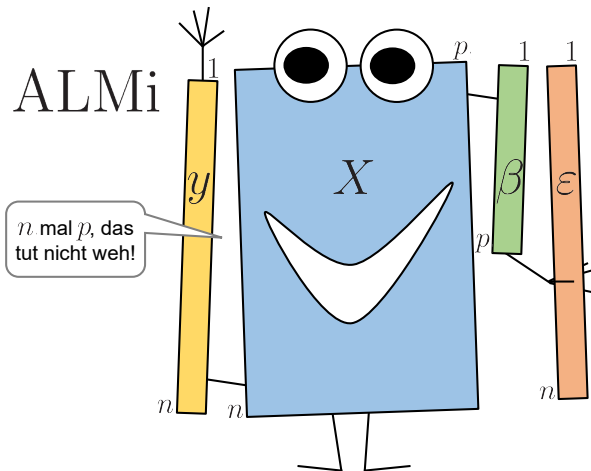
$$y \sim N(I_n 0_n + X\beta, I_n (\sigma^2 I_n) I_n^T) = N(X\beta, \sigma^2 I_n) = N(\mu, \sigma^2 I_n) \quad \text{mit} \quad \mu := X\beta \in \mathbb{R}^n. \quad (13)$$

□

Bemerkungen

- Im ALM sind die Daten y also ein n -dimensionaler normalverteilter Zufallsvektor mit Erwartungswertparameter $\mu = X\beta \in \mathbb{R}^n$ und Kovarianzmatrixparameter $\sigma^2 I_n \in \mathbb{R}^{n \times n}$.
- Die Komponenten y_1, \dots, y_n von y , also die Datenpunkte, sind damit unabhängige, aber im Allgemeinen nicht identisch verteilte, normalverteilte Zufallsvariablen der Form $y_i \sim N(\mu_i, \sigma^2)$ für $i = 1, \dots, n$.





Allgemeine Theorie

Unabhängige und identisch normalverteilte Zufallsvariablen

Einfache lineare Regression

Selbstkontrollfragen

Unabhängige und identisch normalverteilte Zufallsvariablen

Wir betrachten das Szenario von n unabhängigen und identisch normalverteilten Zufallsvariablen mit Erwartungswertparameter $\mu \in \mathbb{R}$ und Varianzparameter σ^2 ,

$$y_i \sim N(\mu, \sigma^2) \quad \text{für } i = 1, \dots, n. \quad (14)$$

Dann gilt, dass (14) äquivalent ist zu

$$y_i = \mu + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \quad \text{für } i = 1, \dots, n \quad \text{mit unabhängigen } \varepsilon_i. \quad (15)$$

In Matrixschreibweise ist dies wiederum äquivalent zu

$$y \sim N(X\beta, \sigma^2 I_n) \quad \text{mit } X := \mathbf{1}_n \in \mathbb{R}^{n \times 1}, \quad \beta := \mu \in \mathbb{R}^1, \quad \sigma^2 > 0. \quad (16)$$

Bemerkungen

- Wir kennen dieses Modell bereits (siehe Einheit (9) in *Wahrscheinlichkeitstheorie und Frequentistische Inferenz*) und haben es bisher geschrieben als

$$y_1, \dots, y_n \sim N(\mu, \sigma^2) \quad \text{mit } (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_{>0}. \quad (17)$$

Unabhängige und identisch normalverteilte Zufallsvariablen

```
# Modellformulierung
library(MASS) # multivariate Normalverteilung
set.seed(1) # reproduzierbare Ergebnisse
n = 12 # Anzahl von Datenpunkten
p = 1 # Anzahl von Betaparametern
X = matrix(rep(1,n), nrow = n) # n x p Designmatrix
I_n = diag(n) # n x n Einheitsmatrix
beta = 2 # wahrer, aber unbekannter Betaparameter
sigsqr = 1 # wahrer, aber unbekannter Varianzparameter

# Datenrealisierung
y = mvrnorm(1, X %*% beta, sigsqr*I_n) # eine Realisierung des n-dimensionalen ZVs y
print(y)
```

```
> [1] 2.39 3.51 1.69 2.58 2.74 2.49 1.18 2.33 3.60 1.16 2.18 1.37
```

Unabhängige und identisch normalverteilte Zufallsvariablen

Designmatrix des Modells u.i.n.v. Zufallsvariablen ($n = 12$, $p = 1$)

Allgemeine Theorie

Unabhängige und identisch normalverteilte Zufallsvariablen

Einfache lineare Regression

Selbstkontrollfragen

Einfache lineare Regression

Wir betrachten das Modell der einfachen linearen Regression (siehe Einheit (1) in *Allgemeines Lineares Modell*),

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \quad \text{für } i = 1, \dots, n, \quad (18)$$

Wir haben bereits gesehen, dass dieses Modell folgende Datenverteilung besitzt:

$$y_i \sim N(\mu_i, \sigma^2) \quad \text{mit } \mu_i := \beta_0 + \beta_1 x_i \quad \text{für } i = 1, \dots, n. \quad (19)$$

In Matrixschreibweise ist dies wiederum äquivalent zu

$$y \sim N(X\beta, \sigma^2 I_n) \quad \text{mit } X := \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \in \mathbb{R}^{n \times 2}, \quad \beta := \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \in \mathbb{R}^2, \quad \sigma^2 > 0. \quad (20)$$

Einfache lineare Regression

```
# Modellformulierung
library(MASS) # multivariate Normalverteilung
set.seed(1) # reproduzierbare Ergebnisse
n = 10 # Anzahl von Datenpunkten
p = 2 # Anzahl von Betaparametern
x = 1:n # Werte des Prädiktors
X = matrix(c(rep(1,n),x), nrow = n) # n x p Designmatrix
I_n = diag(n) # n x n Einheitsmatrix
beta = matrix(c(0,1), nrow = p) # wahrer, aber unbekannter Betaparameter
sigsqr = 1 # wahrer, aber unbekannter Varianzparameter

# Datenrealisierung
y = mvrnorm(1, X %*% beta, sigsqr*I_n) # eine Realisierung des n-dimensionalen ZVs y
print(y)
```

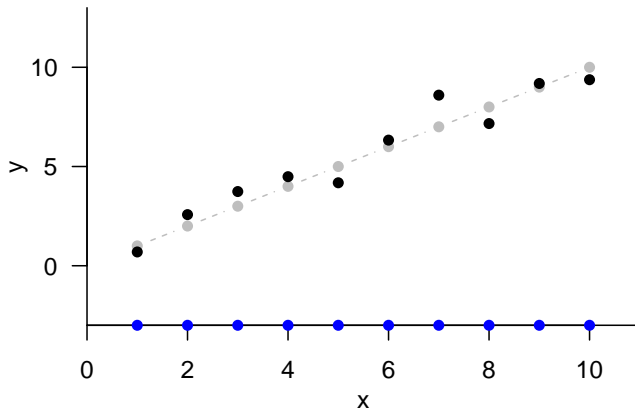
```
> [1] 0.695 2.576 3.738 4.487 4.180 6.330 8.595 7.164 9.184 9.374
```

Einfache lineare Regression

Designmatrix des Modells der einfachen linearen Regression ($n = 10, p = 2$)

Einfache lineare Regression

• x_i • $X\beta$ für $\beta_0 := 0, \beta_1 := 1$ • (x_i, y_i)



Allgemeine Theorie

Unabhängige und identisch normalverteilte Zufallsvariablen

Einfache lineare Regression

Selbstkontrollfragen

Selbstkontrollfragen

1. Erläutern Sie das naturwissenschaftliche Paradigma.
2. Erläutern Sie die Standardprobleme der Frequentistischen Inferenz.
3. Geben Sie die Definition des Allgemeinen Linearen Modells wieder.
4. Erläutern Sie die deterministischen und probabilistischen Aspekte des ALMs.
5. Wieviele skalare Parameter hat das ALM mit sphärischer Kovarianzmatrix?
6. Warum sind die Komponenten des ALM Zufallsfehlers unabhängig und identisch verteilt?
7. Geben Sie das Theorem zur Datenverteilung des Allgemeinen Linearen Modells wieder.
8. Sind die Komponenten des ALM-Datenvektors immer unabhängig und identisch verteilt?
9. Schreiben Sie das Szenario von n unabhängig und identisch normalverteilten Zufallsvariablen in ALM-Form.
10. Schreiben Sie das Szenario der einfachen linearen Regression in ALM-Form.