

# Schätzen und Testen I

Wintersemester 2014/15

Folien zur Vorlesung von

Sonja Greven  
Christian Heumann

nach einer Vorlage von

Christiane Fuchs  
Ludwig Fahrmeir  
Volker Schmid

# Worum geht es in dieser Vorlesung?

Statistik umfasst die

- ▶ deskriptive Statistik
- ▶ explorative Statistik
- ▶ induktive Statistik

Hier: Statistische Inferenz

Verschiedene Inferenzkonzepte. In dieser Vorlesung (ST1) Fokus auf

- ▶ klassischer Inferenz
- ▶ Likelihood-Inferenz
- ▶ Bayes-Inferenz

# Inhalt

1. Einführung in statistische Modelle und Inferenzkonzepte
  - Statistische Modelle
  - Konzepte der statistischen Inferenz
2. Klassische Schätz- und Testtheorie
  - Klassische Schätztheorie
  - Klassische Testtheorie
  - Bereichsschätzung und Konfidenzintervalle
  - Multiples Testen
3. Likelihood-Inferenz
  - Parametrische Likelihood-Inferenz
  - Maximum-Likelihood-Schätzung
  - Testen linearer Hypothesen und Konfidenzintervalle
  - Fehlspezifikation, Quasi-Likelihood und Schätzgleichungen

# Inhalt

1. Einführung in statistische Modelle und Inferenzkonzepte
  - Statistische Modelle
  - Konzepte der statistischen Inferenz
2. Klassische Schätz- und Testtheorie
  - Klassische Schätztheorie
  - Klassische Testtheorie
  - Bereichsschätzung und Konfidenzintervalle
  - Multiples Testen
3. Likelihood-Inferenz
  - Parametrische Likelihood-Inferenz
  - Maximum-Likelihood-Schätzung
  - Testen linearer Hypothesen und Konfidenzintervalle
  - Fehlspezifikation, Quasi-Likelihood und Schätzgleichungen

# 1 Einführung in statistische Modelle und Inferenzkonzepte

## Ziele:

- ▶ Statistische Modelle im Überblick, von einfachen hin zu komplexeren Modellen. Auswahl orientiert an Datenstrukturen, Modellklassen und Fragestellungen aus dem Bachelorprogramm und darüber hinaus.
- ▶ Problemstellungen der zugehörigen statistischen Inferenz.
- ▶ Konzepte statistischer Inferenz im Überblick.

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

Zunächst wird nur der Ein-Stichproben-Fall betrachtet. Seien  $x_1, \dots, x_n$  die Daten als Realisierungen von Stichprobenvariablen  $X_1, \dots, X_n$  und diese i.i.d. wie Zufallsvariable  $X$  mit Verteilungsfunktion  $F(x)$  bzw. (stetiger, diskreter bzw. allgemeiner „Radon-Nikodym“-) Dichte  $f(x)$ .

### Parametrische Modelle

$$X \sim f(x|\theta), \quad \theta = (\theta_1, \dots, \theta_k)^\top \in \Theta \subseteq \mathbb{R}^k$$

In der Regel ist  $k$  fest und klein im Verhältnis zu  $n$ .

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

### Beispiel 1.1

1.  $X \sim N(\mu, \sigma^2)$ ; Schätzen/Testen von  $\mu$ , zum Beispiel Gauß-Test, t-Test; Schätzen/Testen von  $\sigma^2$ .
2. Analoge Problemstellungen für  $X \sim \text{Bin}(n, \pi)$ ,  
 $X \sim \text{Po}(\lambda), \dots$  bzw. allgemein  $X \sim$  einfache lineare Exponentialfamilie mit natürlichem (skalarem) Parameter  $\theta$ .
3.  $\mathbf{X} = (X_1, \dots, X_p)^\top$  mehrdimensional, zum Beispiel  
 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

Beispiel 1.1 fortgeführt

### 4. Lokations- und Skalenmodelle:

$$X \sim F_0 \left( \frac{x - a}{b} \right)$$

mit gegebener Verteilungsfunktion  $F_0(z)$ .

$a \in \mathbb{R}$  heißt Lokationsparameter,  $b > 0$  Skalenparameter.

Dichten im stetigen Fall:

$$X \sim \frac{1}{b} f_0 \left( \frac{x - a}{b} \right)$$

mit gegebener Dichte  $f_0(z)$ .



# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

Beispiele für Lokations- und Skalenmodelle:

- ▶  $X \sim N(a, b^2)$  (Normalverteilung),  $f_0(z) = \phi(z)$ :

$$\frac{1}{b} f_0\left(\frac{x-a}{b}\right) = \frac{1}{\sqrt{2\pi}b} \exp\left(-\frac{1}{2} \frac{(x-a)^2}{b^2}\right)$$

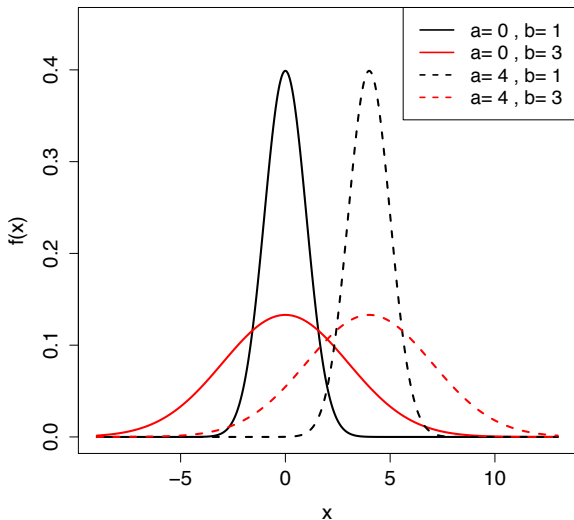
- ▶  $X \sim DE(a, b)$  (Laplace- oder Doppelsexponentialverteilung):

$$\frac{1}{b} f_0\left(\frac{x-a}{b}\right) = \frac{1}{2b} \exp\left(-\frac{|x-a|}{b}\right)$$

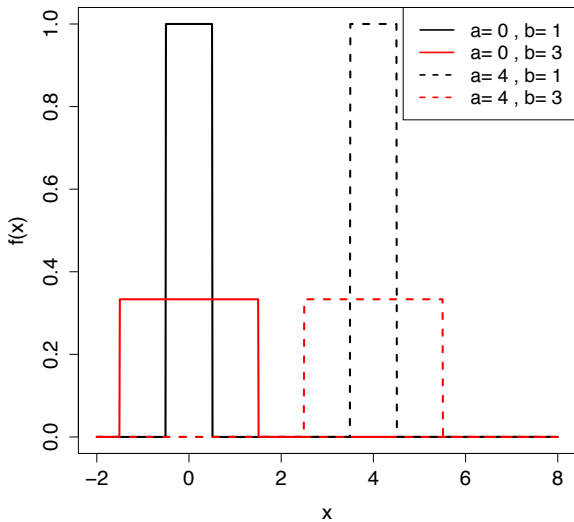
- ▶  $X \sim U(a, b)$  (Gleichverteilung):

$$\frac{1}{b} f_0\left(\frac{x-a}{b}\right) = \frac{1}{b} I_{(a-\frac{b}{2}, a+\frac{b}{2})}(x)$$

Der Träger ist abgeschlossen und hängt von den Parametern ab.



Lokations- und Skalenmodelle: Normalverteilung



Lokations- und Skalenmodelle: Gleichverteilung

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

Weitere Beispiele für Lokations- und Skalenmodelle:

- ▶  $X \sim C(a, b)$  (Cauchy-Verteilung):

$$\frac{1}{b} f_0 \left( \frac{x-a}{b} \right) = \frac{1}{\pi} \cdot \frac{1}{b^2 + (x-a)^2}$$

- ▶  $X \sim L(a, b)$  (logistische Verteilung):

$$\frac{1}{b} f_0 \left( \frac{x-a}{b} \right) = \frac{1}{b} \cdot \frac{\exp \left( -\frac{x-a}{b} \right)}{\left( 1 + \exp \left( -\frac{x-a}{b} \right) \right)^2}$$

- ▶  $X \sim E(a, b)$  (Exponentialverteilung):

$$\frac{1}{b} f_0 \left( \frac{x-a}{b} \right) = \frac{1}{b} \exp \left( -\frac{x-a}{b} \right) I_{[a, \infty)}(x)$$

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

### 5. Exponentialfamilien:

#### Definition 1.1 (Exponentialfamilien)

Eine Verteilungsfamilie heißt Exponentialfamilie  $\stackrel{\text{def}}{\Leftrightarrow}$

$$f(x|\boldsymbol{\theta}) = h(x) \cdot c(\boldsymbol{\theta}) \cdot \exp(\gamma_1(\boldsymbol{\theta})T_1(x) + \dots + \gamma_k(\boldsymbol{\theta})T_k(x)) = \\ h(x) \exp(b(\boldsymbol{\theta}) + \boldsymbol{\gamma}(\boldsymbol{\theta})^\top \mathbf{T}(x))$$

mit  $h(x) \geq 0$  und

$$b(\boldsymbol{\theta}) = \log(c(\boldsymbol{\theta}))$$

$$\mathbf{T}(x) = (T_1(x), \dots, T_k(x))^\top$$

$$\boldsymbol{\gamma}(\boldsymbol{\theta}) = (\gamma_1(\boldsymbol{\theta}), \dots, \gamma_k(\boldsymbol{\theta}))^\top.$$

$\gamma_1, \dots, \gamma_k$  heißen die natürlichen oder kanonischen Parameter der Exponentialfamilie (nach Reparametrisierung von  $\boldsymbol{\theta}$  mit  $\boldsymbol{\gamma}$ ).

Annahme:  $1, \gamma_1, \dots, \gamma_k$  und  $1, T_1(x), \dots, T_k(x)$  sind linear unabhängig, d.h.  $f$  ist strikt  $k$ -parametrisch.

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

### **Beispiel 1.2** (Bernoulli-Experiment)

$$X = (X_1, \dots, X_n) \stackrel{\text{i.i.d.}}{\sim} \text{Bin}(1, \pi).$$

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

**Beispiel 1.2** (Bernoulli-Experiment) fortgeführt

d.h. es liegt eine einparametrische Exponentialfamilie vor mit

$$T(x) = \sum_{i=1}^n x_i$$
$$\gamma = \log\left(\frac{\pi}{1-\pi}\right) =: \text{logit}(\pi).$$

**Bemerkung:** Eine Verteilungsfamilie heißt *einfache lineare Exponentialfamilie*, falls

$$f(x|\theta) \propto \exp(b(\theta) + \theta x)$$

bzw. (mit Dispersionsparameter  $\phi$ ) falls

$$f(x|\theta) \propto \exp\left(\frac{b(\theta) + \theta x}{\phi}\right).$$

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

Beispiel 1.1 fortgeführt

### 6. Mischverteilungen:

$$X \sim \pi_1 f_1(x|\vartheta_1) + \dots + \pi_k f_k(x|\vartheta_k)$$

mit  $\pi_1 + \dots + \pi_k = 1$ , wobei die  $\pi_i$  als *Mischungsanteile* und die  $f_i(x|\vartheta_i)$  als *Mischungskomponenten* bezeichnet werden. Genauer spricht man von *diskreter Mischung*.

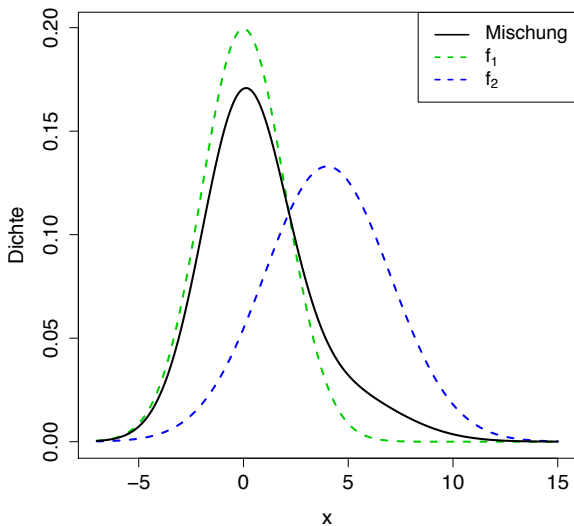
### Beispiel 1.3

$$X \sim \pi_1 \phi(x; \mu_1, \sigma_1^2) + \dots + \pi_k \phi(x; \mu_k, \sigma_k^2)$$

wird *Normalverteilungsmischung* genannt.

Unbekannt sind meistens  $\vartheta = (\vartheta_1, \dots, \vartheta_k)$  und  $\pi = (\pi_1, \dots, \pi_k)$ . Das Schätzen von  $\theta = (\vartheta, \pi)$  erfolgt mit ML-Schätzung, meist mit Hilfe des EM-Algorithmus. Auch gewünscht: Testen auf Anzahl  $k$  der Mischungskomponenten.





Mischung mit  $\pi_1 = 0.8$ ,  $f_1(x) = \phi(x; 0, 2^2)$ ,  $\pi_2 = 0.2$ ,  $f_2(x) = \phi(x; 4, 3^2)$ .

# 1.1 Statistische Modelle

## 1.1.1 Einfache Zufallsstichproben

### **Nichtparametrische Modelle/Inferenz**

- ▶  $X \sim F(x)$ ,  $X$  stetige Zufallsvariable,  $F$  stetige Verteilung
  - ▷ Kolmogorov-Smirnov-Test auf  $H_0 : F(x) = F_0(x)$
- ▶  $X \sim F(x)$ ,  $X$  diskret bzw. gruppiert
  - ▷  $\chi^2$ -Anpassungstest
- ▶  $X \sim f(x)$ ,  $X$  stetige Zufallsvariable,  $f$  bis auf endlich viele Punkte stetig, differenzierbar etc.
  - ▷ nichtparametrische Dichteschätzung, zum Beispiel Kerndichteschätzung

Der Zwei-und Mehr-Stichprobenfall kann analog behandelt werden; vgl. Statistik II.

# 1.1 Statistische Modelle

## 1.1.2 Lineare und generalisierte lineare parametrische Modelle

Daten  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ , sind gegeben, mit  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ .  $y_1 | \mathbf{x}_1, \dots, y_n | \mathbf{x}_n$  sind (bedingt) unabhängig, aber *nicht* identisch verteilt.

### Klassisches lineares Modell (LM)

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} [N](0, \sigma^2) \Leftrightarrow y_i | \mathbf{x}_i \sim [N](\mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$$

- ▶ Annahme:  $p = \dim(\boldsymbol{\beta}) < n$  und  $n$  fest.
- ▶ Schätzen von  $\boldsymbol{\beta}$  und  $\sigma^2$ , Tests für  $\boldsymbol{\beta}$  mit oder ohne Normalverteilungsannahme.
- ▶ Variablenselektion und Modellwahl.  
Spezialfall: Varianzanalyse/Versuchsplanung.

# 1.1 Statistische Modelle

## 1.1.2 Lineare und generalisierte lineare parametrische Modelle

### Generalisierte lineare Modelle (GLM)

$y_i | \mathbf{x}_i$ ,  $i = 1, \dots, n$ , besitzen Dichte aus einfacher linearer Exponentialfamilie, zum Beispiel Normal-, Binomial-, Poisson- oder Gammaverteilung, und sind bedingt unabhängig.

$$\mathbb{E}[y_i | \mathbf{x}_i] = \mu_i = h(\mathbf{x}_i^\top \boldsymbol{\beta})$$

Dabei ist  $h$  die *inverse Linkfunktion* (oder *Responsefunktion*).

# 1.1 Statistische Modelle

## 1.1.2 Lineare und generalisierte lineare parametrische Modelle

### Beispiel 1.4

Sei  $y_i | \mathbf{x}_i \in \{0, 1\}$  und

$$\mu_i = \pi_i = \mathbb{P}(y_i = 1 | \mathbf{x}_i), \quad \pi_i = h(\mathbf{x}_i^\top \boldsymbol{\beta}).$$

Beispiele für  $h$  sind die Verteilungsfunktion der logistischen Verteilung ( $\rightarrow$  Logit-Modell) oder die Verteilungsfunktion der Normalverteilung ( $\rightarrow$  Probit-Modell).

Die Inferenzprobleme im GLM sind wie im linearen Modell. Es ist likelihoodbasierte oder bayesianische Inferenz möglich.

**Beachte:** Die  $y_i | \mathbf{x}_i$  sind nicht identisch verteilt.

# 1.1 Statistische Modelle

## 1.1.3 Nicht- und semiparametrische Regression

### Nichtparametrische Einfachregression

Daten wie im linearen Modell,  $x_i$  skalar.

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$$

Regressionsfunktion  $f(x_i) = \mathbb{E}[y_i|x_i]$  *nicht* parametrisch spezifiziert.

- ▶ Nicht- oder semiparametrisches Schätzen von  $f$
- ▶ Testen von

$$H_0 : f(x) = \beta_0 + x\beta_1 \text{ vs.}$$

$$H_1 : f \text{ nichtlinear.}$$

# 1.1 Statistische Modelle

## 1.1.3 Nicht- und semiparametrische Regression

### Additive Modelle (AM)

$$y_i = f_1(x_{i1}) + \dots + f_p(x_{ip}) + \mathbf{z}_i^\top \boldsymbol{\beta} + \varepsilon_i, \varepsilon_i \text{ wie bisher,}$$

$$\mu_i = \mathbb{E}[y_i | \mathbf{x}_i] = f_1(x_{i1}) + \dots + f_p(x_{ip}) + \mathbf{z}_i^\top \boldsymbol{\beta}.$$

- ▶ Schätzen, Testen von  $f_1, \dots, f_p, \boldsymbol{\beta}$
- ▶ Variablenselektion und Modellwahl (zum Beispiel Einfluss einer bestimmten Kovariable linear oder nichtlinear)

# 1.1 Statistische Modelle

## 1.1.3 Nicht- und semiparametrische Regression

### **Generalisierte Additive Modelle (GAM)**

$y_i | \mathbf{x}_i$  wie bei GLM; analog zu additiven Modellen lässt man aber

$$\mu_i = \mathbb{E}[y_i | \mathbf{x}_i] = h \left( f_1(x_{i1}) + \dots + f_p(x_{ip}) + \mathbf{z}_i^\top \boldsymbol{\beta} \right)$$

zu.



# 1.1 Statistische Modelle

## 1.1.4 Quantil-Regression/Robuste Regression

Datenlage wie bei üblicher Regression:  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ ,  
 $y_i | \mathbf{x}_i$  bedingt unabhängig.

### Ziel:

Schätze nicht (nur)  $\mathbb{E}[y_i | \mathbf{x}_i]$ , zum Beispiel durch KQ-Schätzer  $\mathbf{x}_i^\top \hat{\beta}_{\text{KQ}}$ , sondern den bedingten Median ( $\tau = 0.5$ ) oder allgemeiner die (bedingten) Quantile  $Q_\tau(y_i | \mathbf{x}_i)$ . Statt KQ-Ansatz (ohne Normalverteilungsannahme) und Schätzung von  $\hat{\beta}_{\text{KQ}}$ , so dass  $\sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2$  minimiert wird, suchen wir

$$\hat{\beta}_{\text{med}} := \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n |y_i - \mathbf{x}_i^\top \beta|$$

$$\Rightarrow \mathbf{x}^\top \hat{\beta}_{\text{med}} = \widehat{\operatorname{med}}(y | \mathbf{x}).$$

# 1.1 Statistische Modelle

## 1.1.4 Quantil-Regression/Robuste Regression

Gesucht:

$$\hat{\beta}_{\text{med}} := \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n |y_i - \mathbf{x}_i^{\top} \beta|$$

$$\Rightarrow \mathbf{x}^{\top} \hat{\beta}_{\text{med}} = \widehat{\operatorname{med}}(y|\mathbf{x}).$$

Wichtig dabei: keine Annahme für die Fehlerverteilung, d.h. „verteilungsfreier Ansatz“.

Frage: Welche Konzepte zum Schätzen und Testen verwenden?  
→ Quasi-Likelihood-Methoden.

# 1.1 Statistische Modelle

## 1.1.5 Verweildaueranalyse: Cox-Modell

Grundlegender Begriff:

Hazardrate  $\lambda(t)$  einer stetigen Lebensdauer  $T \geq 0$ .

### Definition 1.2 (Hazardrate)

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{\mathbb{P}(t \leq T \leq t + \Delta t | T \geq t)}{\Delta t}$$
$$\Leftrightarrow \mathbb{P}(t \leq T \leq t + \Delta t | T \geq t) = \lambda(t)\Delta t + o(\Delta t)$$

(Dabei ist  $f(x) = o(g(x))$  für  $x \rightarrow 0$  falls  $\lim_{x \rightarrow 0} f(x)/g(x) = 0$ .)

Interpretation:  $\lambda(t)\Delta t \approx$  bedingte Wahrscheinlichkeit für Ausfall in  $[t, t + \Delta t]$  gegeben „Überleben“ bis zum Zeitpunkt  $t$  bei „kleinem“  $\Delta t$ .  
Mit Kovariablen  $\mathbf{x} = (x_1, \dots, x_p)^\top$ :

$$\lambda(t; \mathbf{x}) = \lim_{\Delta t \rightarrow 0} \frac{\mathbb{P}(t \leq T \leq t + \Delta t | T \geq t; \mathbf{x})}{\Delta t}.$$

# 1.1 Statistische Modelle

## 1.1.5 Verweildaueranalyse: Cox-Modell

### Rechtszensierte Survivaldaten

Verwende  $t_1, \dots, t_n$  für evtl. rechtszensierte Beobachtungen von unabhängigen Lebensdauern  $T_1, \dots, T_n$ ,  $\delta_1, \dots, \delta_n$  als Zensierungsindikatoren und  $\mathbf{x}_1, \dots, \mathbf{x}_n$  als zugehörige Kovariablen.

Ziel: Schätze  $\lambda(t; \mathbf{x})$  bzw. zumindest den Einfluss der Kovariablen auf die Hazardrate.

# 1.1 Statistische Modelle

## 1.1.5 Verweildaueranalyse: Cox-Modell

### Cox-Modell

Im *Cox-Modell* (auch: *Proportional Hazards-Modell*) gilt

$$\begin{aligned}\lambda(t; \mathbf{x}_i) &= \lambda_0(t) \cdot \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) \\ &= \lambda_0(t) \cdot \exp(x_{i1}\beta_1 + \dots + x_{ip}\beta_p) \\ &= \lambda_0(t) \cdot \exp(x_{i1}\beta_1) \cdot \dots \cdot \exp(x_{ip}\beta_p).\end{aligned}$$

Dabei ist  $\lambda_0(t)$  die von  $i$  bzw.  $\mathbf{x}_i$  unabhängige „Baseline“-Hazardrate.  $\exp(\mathbf{x}_i^\top \boldsymbol{\beta})$  modifiziert  $\lambda_0(t)$  multiplikativ.

# 1.1 Statistische Modelle

## 1.1.5 Verweildaueranalyse: Cox-Modell

### **Cox-Modell** (fortgeführt)

Primäres Interesse: Schätzen/Testen von  $\beta$  wie im LM oder GLM;  $\lambda_0(t)$  wird als Nuisanceparameter (bzw. -funktion) betrachtet.

⇒ Die Likelihood faktorisiert sich in

$$L(\beta; \lambda_0(t)) = L_1(\beta) \cdot L_2(\beta; \lambda_0(t)).$$

$L_1(\beta)$  ist *partielle* („partial“) *Likelihood*, die bezüglich  $\beta$  maximiert wird. Erstaunlicherweise ist der Informationsverlust gering. Ferner gibt es einen Zusammenhang zwischen Partial-Likelihood und dem Konzept der Profil-Likelihood.

# 1.1 Statistische Modelle

## 1.1.6 Fehlende/unvollständige Daten

- ▶ Daten: „beliebig“ (Querschnitts-, Survival-, Längsschnittdaten)
- ▶ Beispiele:
  - ▶ Nicht-Antwörter bei statistischen Befragungen
  - ▶ „Drop-out“ bei klinischen Studien
  - ▶ zensierte Daten (wie in Survivalanalyse)
  - ▶ Modelle mit latenten Variablen
- ▶ Übliche Modelle und statistische Methodik setzen vollständige Daten voraus.

# 1.1 Statistische Modelle

## 1.1.7 Konditionale (autoregressive, Markov-) Modelle für Longitudinaldaten

- ▶ **Longitudinaldaten:**  $(y_{ij}, \mathbf{x}_{ij})$  für  $i = 1, \dots, m$  und  $j = 1, \dots, n_i$  als Beobachtungen von Zielvariablen  $y_{ij}$  und Kovariablen  $\mathbf{x}_{ij}$  zu Zeitpunkten  $t_{i1} < \dots < t_{ij} < \dots < t_{in_i}$ . Spezialfall  $m = 1$ : Zeitreihen.
- ▶ **Autoregressives Modell 1. Ordnung bzw. Markov-Modell 1. Ordnung:** Bedingte Verteilung von  $y_{ij} | y_{i,j-1}, y_{i,j-2}, \dots, y_{i1}, \mathbf{x}_{ij}$  ist  $y_{ij} | y_{i,j-1}, \mathbf{x}_{ij}$ , zum Beispiel

$$y_{ij} = \alpha y_{i,j-1} + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \underbrace{\varepsilon_{ij}}_{\text{i.i.d.}}$$

Likelihood-Inferenz: algorithmisch simpel, asymptotische Theorie schwieriger (da  $y_{ij}$  abhängig).



# 1.1 Statistische Modelle

## 1.1.8 (Generalisierte) Lineare gemischte Modelle für Longitudinaldaten

### Lineares gemischtes Modell (LMM)

$$y_{ij} = \beta_0 + \beta_1 t_{ij} + \mathbf{x}_{ij}^\top \boldsymbol{\beta} + \gamma_{0i} + \gamma_{1i} t_{ij} + \varepsilon_{ij}, \quad i = 1, \dots, m, \quad j = 1, \dots, n;$$

- ▶  $\beta_0, \beta_1, \boldsymbol{\beta}$ : feste Populationseffekte, z.B.  $\beta_0 + \beta_1 t$  fester (linearer) Populationstrend
- ▶  $\gamma_{0i}, \gamma_{1i}$ : individuenspezifische Effekte  
⇒ Anzahl der Parameter von der Ordnung des Stichprobenumfangs
- ▶ Annahme:

$$\gamma_{0i} \stackrel{\text{i.i.d.}}{\sim} N(0, \tau_0^2),$$

$$\gamma_{1i} \stackrel{\text{i.i.d.}}{\sim} N(0, \tau_1^2)$$

d.h. die  $\gamma$ -Parameter sind „zufällige“ Parameter.

- ▶ Inferenz: algorithmisch/methodisch variierte Likelihood-Inferenz oder Bayes-Inferenz mit MCMC-Simulationsmethoden. Für GLMM deutlich komplexer als für LMM.

# 1.1 Statistische Modelle

## 1.1.9 Marginale Modelle

→ Kapitel 6.2 und 6.4 bzw. kurze Einführung in 3.4  
(Quasi-Likelihood-Inferenz/GEEs)

# 1.1 Statistische Modelle

## 1.1.10 Modellbasierte Clusteranalyse

- ▶ Idee:  $\mathbf{x} = (x_1, \dots, x_p)^\top$  stammt aus multivariater Mischverteilung mit  $g$  Komponenten:

$$f(\mathbf{x}) = \sum_{k=1}^g p(k) f(\mathbf{x}|\theta_k),$$

zum Beispiel  $f$  Dichte der multivariaten Normalverteilung.

- ▶ Gesucht:
  1. Schätzungen für  $\theta_k, p(k), k = 1, \dots, g$ .
  2. Schätzungen für unbekannte Klassenzugehörigkeit  $k$  eines Objekts mit beobachtetem Merkmalsvektor  $\mathbf{x}$ . Anwendung der Formel von Bayes liefert:

$$\hat{p}(k|\mathbf{x}) = \frac{\hat{p}(k) f(\mathbf{x}|\hat{\theta}_k)}{\hat{f}(\mathbf{x})}.$$

- ▶ Likelihood-Maximierung: mit EM-Algorithmus
- ▶ Bayes: mit MCMC-Algorithmus

# 1.1 Statistische Modelle

## 1.1.11 Modelle mit latenten Variablen

Beobachtet werden Werte  $\mathbf{y}_i = (y_{i1}, \dots, y_{ij}, \dots, y_{ip})^\top$  von  $p$  (korrelierten) Variablen, die als Indikatoren für eine latente, unbeobachtete Variable  $l_i$  (oder eine kleine Zahl von latenten Variablen) dienen.

Primäres Ziel ist die Schätzung der Effekte („Ladungsfaktoren“)  $\lambda_j$  von  $l$  auf den Vektor  $\mathbf{y}$  der Indikatoren, die Schätzung der latenten Werte  $l_i$ ,  $i = 1, \dots, n$ , und die Schätzung der festen Effekte  $\beta$  und  $\gamma$  (siehe nächste Folie).

# 1.1 Statistische Modelle

## 1.1.11 Modelle mit latenten Variablen

### 1. Beobachtungsmodell:

$$y_{ij} = \mathbf{x}_{ij}^{\top} \boldsymbol{\beta} + \lambda_j l_i + \varepsilon_{ij} \quad \text{mit} \quad \varepsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \quad j = 1, \dots, p$$

### 2. Strukturmodell:

$$l_i = \mathbf{u}_i^{\top} \boldsymbol{\gamma} + \delta_i \quad \text{mit} \quad \delta_i \stackrel{i.i.d.}{\sim} N(0, 1)$$

Ohne Kovariablen  $\mathbf{x}$  und  $\mathbf{u}$  ergibt sich das klassische Modell der Faktorenanalyse. Erweiterungen entstehen zum Beispiel durch kategoriale Indikatoren oder nichtlineare Effekte von Kovariablen.

## 1.2 Konzepte der statistischen Inferenz

- ▶  $x = (x_1, \dots, x_n)^\top$  oder  $y = (y_1, \dots, y_n)^\top$  sind Realisierungen von Stichprobenvariablen (Zufallsvariablen)  
 $X = (X_1, \dots, X_n)^\top$  oder  $Y = (Y_1, \dots, Y_n)^\top$ .  
Die Komponenten  $X_1, \dots, X_n$  können auch selbst wieder mehrdimensional sein.
- ▶ Weitere Annahmen:
  - ▶  $X_1, \dots, X_n$  i.i.d. wie  $X \rightarrow$  einfache Zufallsstichprobe (vgl. Abschnitt 1.1.1).
  - ▶  $Y_1, \dots, Y_n$  (bzw.  $Y_1|X_1, \dots, Y_n|X_n$  im Regressionsmodell) sind (bedingt) unabhängig aber *nicht* identisch verteilt.
  - ▶  $Y_1, \dots, Y_n$  sind abhängig, zum Beispiel zeitlich oder räumlich korreliert.

## 1.2 Konzepte der statistischen Inferenz

- ▶ In allen Fällen gilt:  $x \in \mathcal{X}$  bzw.  $y \in \mathcal{Y}$ , wobei  $\mathcal{X}$  bzw.  $\mathcal{Y}$  der entsprechende Stichprobenraum ist.

$X = (X_1, \dots, X_n)^\top$  und  $Y = (Y_1, \dots, Y_n)^\top$  sind auf dem Stichprobenraum nach einer gemeinsamen Verteilung  $\mathbb{P}$  bzw. Verteilungsfunktion  $F(x) = F(x_1, \dots, x_n)$  verteilt.

$\mathbb{P}$  (bzw.  $F$ ) gehört einer Menge (oder Klasse oder Familie) von Verteilungen  $\mathcal{P}_\theta = \{\mathbb{P}_\theta : \theta \in \Theta\}$  an. Zugehörige Verteilungsfunktionen sind  $F(x|\theta)$  bzw. (falls existent) Dichten  $f(x|\theta) = f(x_1, \dots, x_n|\theta)$ .

## 1.2 Konzepte der statistischen Inferenz

Gemeinsame Dichten  $f(x|\theta)$ :

- ▶ i.i.d. Fall:

$$f(x|\theta) = f(x_1|\theta) \cdot \dots \cdot f(x_n|\theta) = \prod_{i=1}^n f(x_i|\theta)$$

- ▶ Unabhängige Zufallsvariablen  $Y_1, \dots, Y_n$ :

$$f(y|\theta) = \prod_{i=1}^n f_i(y_i|\theta),$$

die Dichten hängen also vom Index  $i$  ab.

- ▶ Bei potentiell abhängigen  $Y_1, \dots, Y_n$  ist  $f(y|\theta)$  nicht immer faktorisiert und teils auch analytisch schwer oder nicht darstellbar.



## 1.2 Konzepte der statistischen Inferenz

- ▶ (Übliche) **parametrische** Inferenz:

$$\theta = (\theta_1, \dots, \theta_k)^T \in \Theta \subseteq \mathbb{R}^k, \quad k \text{ fest mit } k < n.$$

- ▶ **Nichtparametrische/verteilungsfreie** Inferenz:

$\Theta$  ist Funktionenraum,  $\theta$  eine bestimmte Funktion. Zum Beispiel ist  $\Theta$  der Raum der stetigen oder differenzierbaren Funktionen.

Beispiele für Methoden: (Kern-)Dichteschätzung, nichtparametrische Regression.

## 1.2 Konzepte der statistischen Inferenz

► **Semiparametrische** Inferenz: Begriff wird nicht ganz einheitlich verwendet für

1.  $\Theta$  hat eine endlich-dimensionale und eine unendlich-dimensionale Komponente. Beispiel: Cox-Proportional-Hazard-Modell.
2. Parameter  $\theta$  hochdimensional, unter Umständen  $\theta = (\theta_1, \dots, \theta_k)^\top$  mit  $k \sim n$ , zum Beispiel bei der semiparametrischen Regression mit Glättungssplines.

Auch:  $k > n$ , zum Beispiel bei GLMs mit Genexpressionsdaten als Kovariablen: Daten  $x_1, \dots, x_k$  mit  $k \sim 1000 - 10000$ , bei nur  $n \sim 50$  Patienten!

Vergleiche multiples Testen in Kapitel 87.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

$X = (X_1, \dots, X_n)$  besitzt Verteilung

$\mathbb{P} \in \mathcal{P} = \{\mathbb{P}_\theta : \theta = (\theta_1, \dots, \theta_k)^\top \in \Theta\}$  mit  $\Theta \subseteq \mathbb{R}^k$  und  $k < n$   
fest, oft  $k \ll n$ .

In der Regel existiert zur Verteilung  $\mathbb{P}_\theta$  eine (diskrete oder stetige bzw. Radon-Nikodym-) Dichte

$$f(x|\theta) = f(x_1, \dots, x_n|\theta).$$

Anmerkung: Allgemein ist dies die Radon-Nikodym-Dichte bezüglich eines dominierenden Maßes, vgl. Maß- und Wahrscheinlichkeitstheorie-Vorlesung.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

► **Punktschätzung:**

Geschätzt werden soll  $\theta$ . Eine messbare Abbildung

$$T : \begin{cases} \mathcal{X} & \longrightarrow \Theta \\ x & \longmapsto T(x) =: \hat{\theta} \end{cases}$$

heißt *Schätzfunktion* oder *Schätzer*.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

► **Punktschätzung** fortgeführt:

Eine Beurteilung der Güte/Optimalität kann z. B. durch

►  $\text{Bias}_\theta(T) = \mathbb{E}_\theta[T] - \theta,$

►  $\text{Var}_\theta(T) = \mathbb{E}_\theta[(T - \mathbb{E}_\theta[T])^2],$

►  $\text{MSE}_\theta(T) = \mathbb{E}_\theta[(T - \theta)^2] = \text{Var}_\theta(T) + (\text{Bias}_\theta(T))^2$

erfolgen.

Das Konzept der „Güte“ ist frequentistisch, da beurteilt wird, wie „gut“  $T = T(X)$  bei „allen“ denkbaren wiederholten Stichproben  $x$  als Realisierung von  $X$  „im Schnitt“ funktioniert. Anders ausgedrückt: Beurteilt wird nicht die konkret vorliegende Stichprobe, sondern (in der Häufigkeitsinterpretation) das „Verfahren“  $T = T(X)$ .

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

### ► Bereichsschätzung / Intervallschätzung:

$$C : \begin{cases} \mathcal{X} & \longrightarrow \mathcal{P}(\Theta) \\ x & \longmapsto C(x) \subseteq \Theta \end{cases}$$

so dass  $\mathbb{P}_\theta(C(X) \ni \theta) \geq 1 - \alpha$  für alle  $\theta \in \Theta$ .

Dabei ist  $1 - \alpha$  der *Vertrauensgrad* (auch: Konfidenzniveau oder Überdeckungswahrscheinlichkeit) des *Konfidenzbereiches*.

Man beachte die frequentistische/Häufigkeitsinterpretation:  $C(X)$  ist ein *zufälliger* Bereich.

Ist  $\Theta \subseteq \mathbb{R}$  und  $C(x)$  für alle  $x$  ein Intervall, dann heißt  $C$  *Konfidenzintervall*.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

- ▶ **Testen:** Mit einem Test  $\phi$  soll eine Hypothese  $H_0$  gegen eine Alternativhypothese  $H_1$  geprüft werden:

$$H_0 : \theta \in \Theta_0 \quad \text{vs.} \quad H_1 : \theta \in \Theta_1,$$

wobei  $\Theta_0 \cap \Theta_1 = \emptyset$ .

Es muss nicht notwendigerweise  $\Theta = \Theta_0 \cup \Theta_1$  gelten.

Ergebnisse/Aktionen:

$A_0$  :  $H_0$  wird nicht abgelehnt,

$A_1$  :  $H_1$  wird bestätigt, das Ergebnis „ist signifikant“.



# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

- ▶ **Testen** fortgeführt:

Der Test ist eine Abbildung

$$\phi : \mathcal{X} \rightarrow \{A_0, A_1\} = \{0, 1\}.$$

Ein nicht-randomisierter Test hat die Form

$$\phi(x) = \begin{cases} 1, & \text{falls } x \in K, \\ 0, & \text{falls } x \notin K. \end{cases}$$

Dabei ist  $K \subset \mathcal{X}$  der sogenannte *kritische Bereich* und als eine Teilmenge aller möglichen Stichproben zu verstehen. Oft formuliert man dies über eine Teststatistik  $T(x)$ :

$$\phi(x) = \begin{cases} 1, & \text{falls } T(x) \in C, \\ 0, & \text{falls } T(x) \notin C. \end{cases}$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

- ▶ **Testen** fortgeführt:

*Test zum Niveau („size“)  $\alpha$ , wobei  $\alpha$  „klein“:*

$$\mathbb{P}_\theta(A_1) \leq \alpha \text{ für alle } \theta \in \Theta_0.$$

Dabei ist die Wahrscheinlichkeit für den *Fehler 1. Art* kleiner als  $\alpha$ . Die Funktion

$$g_\phi(\theta) = \mathbb{P}_\theta(A_1) = \mathbb{E}_\theta[\phi(X)]$$

heißt *Gütefunktion* von  $\phi$ . Synonym zum Begriff Güte werden auch die Begriffe *Power* oder *Macht* gebraucht. Die Forderung für den Fehler formuliert über die Gütefunktion lautet

$$g_\phi(\theta) \leq \alpha \text{ für } \theta \in \Theta_0.$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

► **Testen** fortgeführt:

„Programm“ der klassischen parametrischen Schätztheorie (siehe Kapitel 2):

Finde Test  $\phi$  zum Niveau  $\alpha$  mit „optimaler“ Power bzw. minimaler Wahrscheinlichkeit für den *Fehler 2. Art*,

$1 - g_\phi(\theta)$ ,  $\theta \in \Theta_1$ . Das Konzept ist wiederum frequentistisch.

Das „Programm“ ist dabei hauptsächlich für spezielle Verteilungsfamilien (zum Beispiel für Exponentialfamilien) und spezielle Testprobleme im i.i.d. Fall durchführbar. Weniger tauglich ist es für (etwas) komplexere Modelle, zum Beispiel bereits für GLMs. Dann:

- Likelihood-Inferenz
- Bayes-Inferenz
- Nicht- und semiparametrische Inferenz

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

- ▶ **Testen** fortgeführt:

Im einfachsten Fall von zwei Punkthypothesen

$$H_0 : \theta = \theta_0, \quad H_1 : \theta = \theta_1$$

für  $\theta_0 \neq \theta_1$  hat der „beste“ Test

Likelihood-Quotienten-Struktur:  $H_0$  wird abgelehnt, falls

$$\frac{f(x|\theta_1)}{f(x|\theta_0)} > k_\alpha$$

(vgl. Neyman-Pearson Theorem, Abschnitt 2.2 oder Wahrscheinlichkeitstheorie und Inferenz II).

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

### ► p-Werte (p-values):

#### Beispiel 1.5 (Gauß-Test)

$X_1, \dots, X_n$  i.i.d.  $N(\mu, \sigma^2)$ ,  $\sigma^2$  bekannt. Betrachte

$$H_0 : \mu \leq \mu_0 \quad , \quad H_1 : \mu > \mu_0.$$

Teststatistik ist

$$T(X) = \frac{\bar{X} - \mu_0}{\sigma} \sqrt{n} \stackrel{\mu = \mu_0}{\sim} N(0, 1).$$

$H_0$  wird abgelehnt, wenn  $T(x) > z_{1-\alpha}$ . Der p-Wert ist  
 $p = \mathbb{P}(T(X) > T(x) | \mu = \mu_0) = \sup_{\mu} \mathbb{P}(T(X) > T(x) | H_0)$ .

Offensichtlich gilt:

$$T(x) > z_{1-\alpha} \Leftrightarrow p < \alpha.$$

Der p-Wert liefert mehr Information (nämlich wie nahe man an der Entscheidungsgrenze ist) als die reine „Bekanntgabe“ der Entscheidung.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.1 Klassische parametrische Inferenz

### Definition 1.3 (p-Wert)

Gegeben sei ein Test bzw. eine Teststatistik  $T(X)$  für  $H_0$  vs.  $H_1$  mit

1.  $\sup_{\theta \in \Theta} \mathbb{P}_{\theta}(T(X) \in C_{\alpha} | H_0) \leq \alpha$ ,
2. für  $\alpha \leq \alpha'$  gilt  $C_{\alpha} \subseteq C_{\alpha'}$ .

Dann gilt  $p = \inf\{\alpha : T(x) = t \in C_{\alpha}\}$ , und  $H_0$  wird abgelehnt, falls  $p < \alpha$ .

# 1.2 Konzepte der statistischen Inferenz

## 1.2.2 (Parametrische) Likelihood-Inferenz

- ▶ Sei  $\mathcal{P} = \{f(x|\theta) | \theta \in \Theta\}$ , d.h. es existieren Dichten zu der vorgegebenen parametrisierten Verteilungsfamilie  $\mathcal{P}$ . Nach der Beobachtung von  $X = x$  heißt

$$L(\theta|x) := f(x|\theta)$$

*Likelihoodfunktion* von  $\theta$  zur Beobachtung  $x$ .

# 1.2 Konzepte der statistischen Inferenz

## 1.2.2 (Parametrische) Likelihood-Inferenz

- ▶ Likelihoodprinzip: Besitzen zwei Beobachtungen  $x$  und  $\tilde{x}$  zueinander proportionale Likelihoodfunktionen, sollen sie zu denselben Schlüssen über  $\theta$  führen.

**Beispiel:**  $X_1, \dots, X_n$  i.i.d.  $N(\mu, \sigma^2)$ ,  $\sigma^2$  bekannt.

$$f(x|\mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Zwei Beobachtungen  $x$  und  $y$  mit  $\bar{x} = \bar{y}$  führen nach dem Likelihood-Prinzip zu den gleichen Schlüssen über  $\mu$ .



# 1.2 Konzepte der statistischen Inferenz

## 1.2.2 (Parametrische) Likelihood-Inferenz

- ▶ Punktschätzung: Maximum-Likelihood- (ML-) Schätzung

$$T(x) = \hat{\theta}_{\text{ML}} \text{ mit } f(x|\hat{\theta}_{\text{ML}}) = \max_{\theta} f(x|\theta)$$

bzw.

$$\hat{\theta}_{\text{ML}} = \operatorname{argmax}_{\theta} f(x|\theta).$$

- ▶ In der Regel existieren keine finiten Optimalitätseigenschaften, jedoch asymptotische.
- ▶ Testen: Likelihood-Quotienten-Test, Wald-Test, Score-Test.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.3 Likelihoodbasierte Inferenz

Quasi-Likelihood-Inferenz, penalisierte Likelihood-Inferenz, semiparametrische Modelle.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.4 Bayes-Inferenz

Wir betrachten wieder  $\mathcal{P} = \{f(\mathbf{x}|\theta) : \theta \in \Theta\}$ , zusätzlich wird aber die Unsicherheit über  $\theta$  durch die *Prioridichte*  $p(\theta)$  auf  $\Theta$  bewertet. Dabei kann  $\Theta$  auch sehr hochdimensional sein.

- ▶ Prinzip: Nach Beobachtung von  $\mathbf{x}$  ist sämtliche Information über  $\theta$  enthalten in der *Posterioridichte*

$$p(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta) \cdot p(\theta)}{\int f(\mathbf{x}|\theta) \cdot p(\theta) d\theta} \quad \begin{array}{l} \text{proportional bzgl.} \\ \text{Parameter } \theta \\ \propto \end{array} \quad f(\mathbf{x}|\theta) \cdot p(\theta)$$
$$= L(\theta|\mathbf{x}) \cdot p(\theta).$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.4 Bayes-Inferenz

### ► Bayes-Schätzung:

- Posteriori-Erwartungswert:

$$T_{\mathbb{E}}(x) = \hat{\theta}_{\text{post-EW}} = \mathbb{E}_{\theta|\mathbf{x}}(\theta|\mathbf{x}) = \int_{\Theta} \theta p(\theta|\mathbf{x}) d\theta$$

- Posteriori-Median:

$$T_{\text{med}}(x) = \hat{\theta}_{\text{post-Med}} = \text{med}_{\theta|\mathbf{x}}(\theta|\mathbf{x})$$

- Posteriori-Modus:

$$T_{\text{mod}}(x) = \hat{\theta}_{\text{post-Mod}} = \underset{\theta}{\operatorname{argmax}} p(\theta|\mathbf{x}) = \underset{\theta}{\operatorname{argmax}} p(\theta)L(\theta|\mathbf{x})$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.4 Bayes-Inferenz

- ▶ Es sind auch *uneigentliche Prioriverteilungen* zulässig, d.h. Dichten mit

$$\int_{\Theta} p(\theta) d\theta = +\infty,$$

die sich somit nicht normieren lassen. Allerdings muss die Posterioridichte eigentlich sein!

Ein Spezialfall ist  $p(\theta) \propto 1$  („Gleichverteilungs-Priori“), bei deren Verwendung

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} L(\theta|\mathbf{x}) = \hat{\theta}_{\text{post-Mod}}$$

gilt, d.h. der ML-Schätzwert und der Posteriori-Modus-Schätzwert identisch sind.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.4 Bayes-Inferenz

- ▶ Die Verteilung zu  $p(\theta)$  heißt die *konjugierte Verteilung* für  $f(\mathbf{x}|\theta)$ , wenn  $f(\theta|\mathbf{x})$  (posteriori) und  $f(\theta)$  (priori) dieselbe Form haben, d.h. wenn Priori- und Posteriorverteilung zur selben Verteilungsfamilie gehören.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.4 Bayes-Inferenz

- ▶ Bayes-Bereichsschätzung: Wähle *Kredibilitätsbereiche/-intervalle*  $C(\mathbf{x})$  so, dass

$$\int_{C(\mathbf{x})} p(\theta|\mathbf{x}) d\theta = \mathbb{P}_{\theta|\mathbf{x}} \left( \underbrace{\theta}_{\text{zufällig}} \in \underbrace{C(\mathbf{x})}_{\text{nicht zufällig, deterministisch}} \right) \geq 1 - \alpha.$$

Es ist also eine Wahrscheinlichkeitsaussage für eine konkrete Stichprobe möglich und keine Häufigkeitsinterpretation notwendig!

- ▶ Bei Bayes-Inferenz wird keine Häufigkeitsinterpretation *benötigt*. Allerdings kann sie trotzdem gemacht werden. (→ Asymptotik der Bayes-Schätzer)

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Sichtweise in der Entscheidungstheorie: Schätzen und Testen als Entscheidung unter Unsicherheit.

Wie bisher betrachten wir  $\mathbb{P} \in \mathcal{P}_\theta = \{\mathbb{P}_\theta : \theta = (\theta_1, \dots, \theta_k)^\top \in \Theta\}$  als statistisches Modell;  $x$  bezeichne eine Stichprobe / konkrete Beobachtung von  $X$ .



# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Zusätzlich werden folgende Funktionen betrachtet:

### Definition 1.4 (Entscheidungsfunktion)

*Als Entscheidungsfunktion bezeichnet man eine Funktion*

$$d : \begin{cases} \mathcal{X} & \longrightarrow & D \\ x & \longmapsto & d(x). \end{cases}$$

*Mit  $D$  wird der Entscheidungs- oder Aktionenraum bezeichnet.*

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Definition 1.5 (Verlustfunktion)

*Eine Verlustfunktion (oft auch stattdessen Gewinnfunktion)*

$$L : \begin{cases} D \times \Theta & \longrightarrow \mathbb{R} \\ (d, \theta) & \longmapsto L(d, \theta) \end{cases}$$

*ordnet einer Entscheidung  $d(x)$  („decision“) einen Verlust („loss“) zu. Im Allgemeinen ist  $L$  so gewählt, dass der Verlust bei richtiger Entscheidung null ist, also  $L$  eine nicht-negative Funktion ist.*

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Beispiel 1.6

1. **Test:** Betrachte

$$H_0 : \theta \leq \theta_0 \quad \text{vs.} \quad H_1 : \theta > \theta_0$$

(zum Beispiel Gauß-Test).

Der Entscheidungsraum sei  $D = \{d_0, d_1\}$  mit

$d_0$ : Entscheidung für  $H_0$ ,

$d_1$ : Entscheidung für  $H_1$ .

Eine mögliche Verlustfunktion ist:

$$L(d_0, \theta) = \begin{cases} 0, & \text{falls } \theta \leq \theta_0 & \text{(Entscheidung richtig)} \\ a \in \mathbb{R}_+, & \text{falls } \theta > \theta_0 & \text{(Fehler 2. Art)} \end{cases}$$
$$L(d_1, \theta) = \begin{cases} 0, & \text{falls } \theta > \theta_0 & \text{(Entscheidung richtig)} \\ b \in \mathbb{R}_+, & \text{falls } \theta \leq \theta_0 & \text{(Fehler 1. Art)} \end{cases}$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

**Beispiel 1.6** fortgeführt

2. **Schätzung:** „Entscheidung“ ist reelle Zahl:

$$d(x) = T(x) = \hat{\theta} \in \Theta, \text{ d.h. } D = \Theta.$$

Mögliche Verlustfunktionen:

$$L(d, \theta) = (d - \theta)^2 \quad \text{quadratischer Verlust,}$$

$$L(d, \theta) = |d - \theta| \quad \text{absoluter Verlust,}$$

$$L(d, \theta) = w(\theta)(d - \theta)^2 \quad \text{gewichteter quadratischer Verlust,}$$

wobei  $w$  eine feste Gewichtsfunktion ist.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

**Beispiel 1.6** fortgeführt

3. **Mehrentscheidungsverfahren**, zum Beispiel Wahl zwischen drei Alternativen

$$d_0 : \theta \leq \theta_0, \quad d_1 : \theta > \theta_1, \quad d_2 : \theta_0 < \theta \leq \theta_1.$$

4. Analog: Modellwahl, Variablenselektion

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Definition 1.6 (Risikofunktion)

*Eine Risikofunktion ist definiert als*

$$R(d, \theta) = \mathbb{E}_\theta[L(d(X), \theta)] = \int_{\mathcal{X}} L(d(x), \theta) f(x|\theta) dx$$

*(„Verlust im Mittel“). Sie ist unabhängig von  $x$ . Dabei wird  $d(X)$  rausintegriert, d.h.  $R(d; \theta)$  ist bei gegebenem  $d$  nur noch eine Funktion von  $\theta$ .*

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Beispiel 1.7

#### 1. Schätzen, d.h.

$d(x) = T(x)$  Schätzwert,  $d(X) = T(X)$  Punktschätzer.

Bei quadratischer Verlustfunktion ist

$$L(T(X), \theta) = (T(X) - \theta)^2$$

mit Risikofunktion

$$R(T, \theta) = \mathbb{E}_\theta[(T(X) - \theta)^2] = \text{MSE}_\theta(T(X)).$$

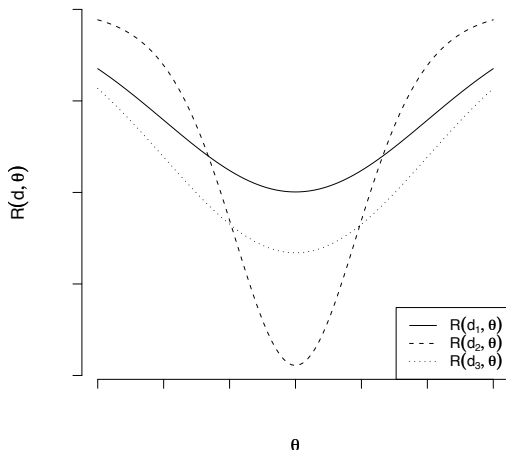
Man beachte, dass das Argument  $T$  in  $R(T, \theta)$  den Schätzer und nicht den konkreten Schätzwert bezeichnet.

#### 2. Testen: vgl. Übung.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Vergleich von Entscheidungsregeln mittels der Risikofunktion



Aus der Abbildung geht hervor, dass  $d_3$  besser als  $d_1$  ist für alle  $\theta \in \Theta$ , d.h.  $d_3$  dominiert  $d_1$  gleichmäßig.

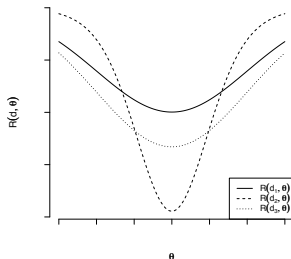


# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

**Ziel:** Finde Regel  $d^*$ , die alle „konkurrierenden“ Regeln  $d$  dominiert.

**Problem:** Diese Idee funktioniert im Allgemeinen nicht, in der Regel überschneiden sich die Risikofunktionen, zum Beispiel ist in obiger Abbildung  $d_2$  nur in einem gewissen Bereich besser als  $d_1$  und  $d_3$ .



# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

→ „Optimale“ Entscheidungsregeln nur möglich durch:

- ▶ Einschränkung auf spezielle Klassen von Verlustfunktionen,
- ▶ Einschränkung auf spezielle Klassen von Entscheidungsregeln, zum Beispiel unverzerrter Schätzer oder unverfälschter Test,
- ▶ oder zusätzliches Kriterium.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Kriterien für "Optimale" Entscheidungsregeln

### 1. **Minimax-Kriterium**

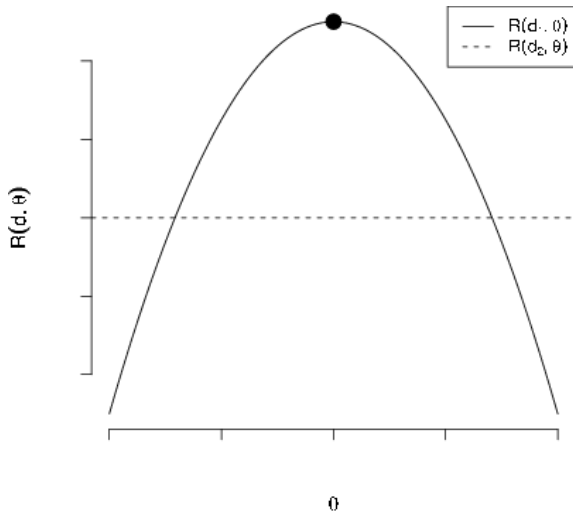
Idee: Betrachte Maximum der Risikofunktion, d.h. präferiere in der folgenden Abbildung  $d_2$ , da

$$\max_{\theta} R(d_2, \theta) < \max_{\theta} R(d_1, \theta).$$

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

$$\max_{\theta} R(d_2, \theta) < \max_{\theta} R(d_1, \theta).$$



# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Definition 1.7 (Minimax-Entscheidungsregel)

Sei  $d^* : \mathcal{X} \rightarrow D$  eine Entscheidungsregel.  $d^*$  heißt *Minimax*, falls es das *supremale Risiko* minimiert:

$$\sup_{\theta \in \Theta} R(d^*, \theta) \leq \sup_{\theta \in \Theta} R(d, \theta) \quad \forall d \in D \Leftrightarrow d^* = \operatorname{arginf}_{d \in D} \sup_{\theta \in \Theta} R(d, \theta).$$

**Bemerkung.** In vielen Fällen werden Supremum und Infimum auch angenommen, so dass tatsächlich

$$d^* = \operatorname{argmin}_{d \in D} \max_{\theta \in \Theta} R(d, \theta)$$

gilt, daher auch der Name Minimax.

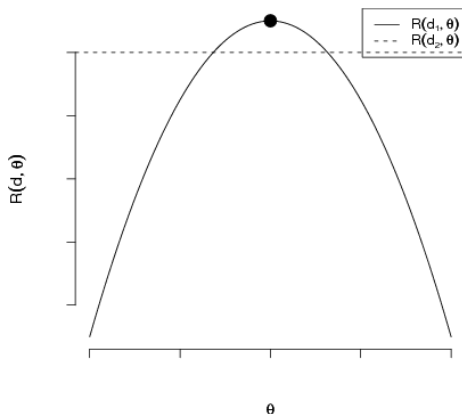
# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Beim Minimax-Kriterium schützt man sich gegen den schlimmsten Fall, was aber nicht unbedingt immer vernünftig ist, wie die folgende Abbildung zeigt:

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie



Hier wäre  $d^*$  nur dann vernünftig, wenn  $\theta$ -Werte in der Mitte „besonders wahrscheinlich“ sind.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Kriterien für "Optimale" Entscheidungsregeln

### 2. Bayes-Kriterium

Wie in der Bayes-Inferenz nehmen wir für  $\theta$  eine Prioridichte  $p(\theta)$  an (aus frequentistischer Sichtweise ist  $p(\theta)$  eine – nicht notwendigerweise normierte – Gewichtsfunktion).

Das *Bayes-Risiko* ist

$$\begin{aligned}r(d, p) &= \int_{\Theta} R(d, \theta) p(\theta) d\theta \\ &= \mathbb{E}_p[R(d, \theta)] \\ &= \mathbb{E}_p \mathbb{E}_{\theta}[L(d(X), \theta)] \\ &= \int_{\Theta} \int_{\mathcal{X}} L(d(x), \theta) f(x|\theta) dx p(\theta) d\theta\end{aligned}$$



# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### 2. **Bayes-Kriterium** fortgeführt:

Das *Bayes-Risiko* wird durch den *Bayes-optimalen Schätzer*  $d^*$  minimiert:

$$r(d^*, p) = \inf_{d \in D} r(d, p).$$

Sei  $p(\theta|x)$  (eigentliche) Posterioridichte. Dann heißt

$$\int_{\Theta} L(d(x), \theta) p(\theta|x) d\theta = \mathbb{E}_{\theta|x}[L(d(x), \theta)]$$

das *Posteriori-Bayes-Risiko*.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Es gilt folgendes praktische Resultat:

### Satz 1.8

*Eine Regel  $d^*$  ist genau dann Bayes-optimal, wenn  $d^*$  für jede Beobachtung/Stichprobe  $x$  das Posteriori-Bayes-Risiko minimiert.*

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

Anmerkungen:

- ▶ Satz 1.8. erleichtert die Bestimmung der Bayes-optimalen Entscheidungsregel in konkreten Fällen.
- ▶ Er zeigt eine intuitive Eigenschaft des Bayesianischen Vorgehens: Um eine Entscheidung geg. eine Beobachtung  $x$  zu treffen, reicht es, den Verlust für  $d(x)$  zu betrachten. Es ist nicht nötig, Verluste  $d(X)$  für andere mögliche aber nicht beobachtete  $X$  zu berücksichtigen.
- ▶ Bayes-optimale Regeln  $d^*$  sind *zulässig*, d.h. sie werden von keiner anderen Regel  $d \neq d^*$  dominiert.
- ▶ Eine enge Beziehung zur Minimax-Regel ist durch die Wahl einer „ungünstigsten“ Prioridichte  $p^*(\theta)$  herstellbar.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.5 Statistische Entscheidungstheorie

### Optimalität von Bayes-Schätzern:



$$\hat{\theta} = \mathbb{E}[\theta|x] = \int_{\Theta} \theta p(\theta|x) d\theta$$

ist Bayes-optimal bei quadratischer Verlustfunktion  
 $L(d, \theta) = (d - \theta)^2$ .



$$\hat{\theta} = \text{med}(\theta|x)$$

ist Bayes-optimal bei absoluter Verlustfunktion  
 $L(d, \theta) = |d - \theta|$ .

## 1.2 Konzepte der statistischen Inferenz

### 1.2.5 Statistische Entscheidungstheorie

**Optimalität von Bayes-Schätzern** fortgeführt:



$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} p(\theta|x)$$

ist Bayes-optimal bei 0-1 Verlustfunktion

$$L_{\varepsilon}(d, \theta) = \begin{cases} 1, & \text{falls } |d - \theta| \geq \varepsilon, \\ 0, & \text{falls } |d - \theta| < \varepsilon \end{cases}$$

für  $\varepsilon > 0$  genügend klein.

Der Grenzübergang  $\varepsilon \rightarrow 0$  liefert den Posteriori-Modus.

- ▶ Anmerkung: Die ML-Schätzung ist optimal bei flacher Priori  $p(\theta) \propto 1$  und bei Wahl obiger 0-1-Verlustfunktion.

# 1.2 Konzepte der statistischen Inferenz

## 1.2.6 Weitere Inferenzkonzepte

- ▶ Struktur-Inferenz
- ▶ Fiduzial-Inferenz