

Mathias Trabs  
Moritz Jirak  
Konstantin Krenz  
Markus Reiß

# Statistik und maschinelles Lernen

Eine mathematische Einführung in klassische  
und moderne Methoden

Vorläufige Arbeitsversion vom 23. Dezember 2025

Springer Nature



# Inhaltsverzeichnis

<b>1</b>	<b>Grundlagen der Statistik</b>	1
1.1	Das statistische Modell	1
1.2	Parameterschätzung	5
1.2.1	Konstruktionsprinzipien	5
1.2.2	Minimax- und Bayes-Ansatz	13
1.3	Hypothesentests	20
1.3.1	Statistische Tests und ihre Fehler	20
1.3.2	Das Neyman-Pearson-Lemma	33
1.4	Konfidenzmengen	40
1.5	Aufgaben	44
<b>2</b>	<b>Das lineare Modell</b>	49
2.1	Regression und kleinste Quadrate	49
2.1.1	Lineare Regression	50
2.1.2	Schätzen im linearen Modell	55
2.1.3	Zufälliges Design und Vorhersage	65
2.2	Inferenz unter Normalverteilungsannahme	73
2.3	Varianzanalyse	87
2.4	Aufgaben	95
<b>3</b>	<b>Modellwahl</b>	101
3.1	Informationskriterien	101
3.1.1	Akaike-Informationskriterium (AIC)	106
3.1.2	Das Bayes-Informationskriterium (BIC)	113
3.2	Orakelungleichung für die penalisierte Modellwahl	117
3.3	Aufgaben	122
<b>A</b>	<b>Konzepte der Wahrscheinlichkeitstheorie</b>	127
A.1	Grundbegriffe der Maßtheorie und Stochastik	127
A.2	Diskrete Verteilungen	142
A.3	Stetige Verteilungen	144

<b>Literaturverzeichnis</b> .....	147
<b>Sachverzeichnis</b> .....	149

## Abkürzungen und Symbole

$\emptyset$	leere Menge
$A^C$	Komplement der Menge $A$
$ A $	Kardinalität der endlichen Menge $A$
$\mathbb{1}_A(x), \mathbb{1}(x \in A)$	Indikatorfunktion zur Menge $A$ , das heißt falls $x \in A$ , dann gilt $\mathbb{1}_A(x) = \mathbb{1}(x \in A) = 1$ und sonst 0
$\mathcal{P}(A)$	Potenzmenge der Menge $A$
$\mathcal{B}(\Omega)$	Borel- $\sigma$ -Algebra über Menge $\Omega$
$\mathbb{N}, \mathbb{N}_0$	Menge der natürlichen Zahlen ohne bzw. mit der Null
$\mathbb{R}, \mathbb{R}_+$	Menge der reellen Zahlen bzw. Menge der reellen Zahlen größer oder gleich Null
$a \wedge b, a \vee b$	Minimum bzw. Maximum von $a, b \in \mathbb{R}$
$\arg \min_a f(a)$	Minimalstelle oder Menge der Minimalstellen der Funktion $f$
$\arg \max_a f(a)$	Maximalstelle oder Menge der Maximalstellen der Funktion $f$
$(a)_+, (a)_-$	Positiv- bzw. Negativteil, das heißt $(a)_+ = a \vee 0, (a)_- = (-a) \vee 0$
$\propto$	proportional, das heißt Gleichheit bis auf eine multiplikative Konstante
$(\mathcal{X}, \mathcal{F}, \mathbb{P}_\vartheta)$	Wahrscheinlichkeitsraum bezüglich des Parameters $\vartheta \in \Theta$
$\mathbb{E}_\vartheta, \text{Var}_\vartheta$	Erwartungswert bzw. Varianz bezüglich $\mathbb{P}_\vartheta$
$\text{Cov}, \mathbb{Cov}$	Kovarianz und Kovarianzmatrix
$U(A)$	Gleichverteilung auf der Menge $A$
$N(\mu, \sigma^2)$	Normalverteilung mit Mittelwert $\mu$ und Varianz $\sigma^2$
$\text{Ber}(\vartheta)$	Bernoulli-Verteilung mit Erfolgswahrscheinlichkeit $\vartheta \in [0, 1]$
$\Phi$	Verteilungsfunktion von $N(0, 1)$
$\sim$	ist verteilt gemäß, zum Beispiel „ $Y \sim N(\mu, \sigma^2)$ “
i.i.d.	unabhängig und identisch verteilt (englisch: <i>independent and identically distributed</i> )

$\bar{X}_n$	Stichprobenmittel $\bar{X}_n = \frac{1}{n} \sum_{j=1}^n X_j$ einer Stichprobe $X_1, \dots, X_n$
$q_\alpha$	$\alpha$ -Quantil der Standardnormalverteilung
$q_{V,\alpha}$	$\alpha$ -Quantil einer Verteilung $V \in \{U(A), N(\mu, \sigma^2), \dots\}$
$ \cdot , \langle \cdot, \cdot \rangle$	Euklidische Norm und euklidisches Skalarprodukt auf $\mathbb{R}^d$
$ \cdot _p$	Vektor- oder Folgenorm für $p \in [0, \infty]$
$\ \cdot\ , \ \cdot\ _2$	Spektralnorm bzw. Frobeniusnorm einer Matrix
$\mathbb{R}^{m \times n}$	Menge der reellen Matrizen mit $m \in \mathbb{N}$ Zeilen und $n \in \mathbb{N}$ Spalten
$E_n$	Einheitsmatrix mit $n \in \mathbb{N}$ Spalten und Zeilen
$\mathcal{L}^p, L^p$	Raum der $p$ -fach integrierbaren Funktionen bzw. Äquivalenzklassen
span	Spann, das heißt der lineare Unterraum aufgespannt durch die gegebenen Vektoren
Im	Bildbereich einer Funktion (englisch: <i>image</i> )
tr	Spur einer Matrix (englisch: <i>trace</i> )
rank	Rang einer Matrix
sgn	Vorzeichen einer reellen Zahl, d. h. $\text{sgn}(x) = 1$ , falls $x > 0$ und sonst $-1$ .

Weitere Symbole und Begriffe aus der Wahrscheinlichkeitstheorie sind im Anhang zu finden.

# Kapitel 1

## Grundlagen der Statistik

In diesem Kapitel werden die grundlegenden Begriffe und Konzepte der Statistik eingeführt. Der Startpunkt ist die mathematische Formulierung eines statistischen Modells. Davon ausgehend werden wir Konstruktionsprinzipien für Parameterschätzer, Hypothesentests und Konfidenzbereiche kennenlernen. Nachdem diese in den einfachen Beispielen dieses Kapitels verstanden sind, können wir aus diesen Grundprinzipien heraus neue Methoden in den komplexeren Modellen späterer Kapitel entwickeln.

### 1.1 Das statistische Modell

Während die Wahrscheinlichkeitstheorie anhand eines gegebenen Modells die Eigenschaften der (zufälligen) Ereignisse untersucht, ist das Ziel der Statistik entgegengesetzt: Wie kann man aus den gegebenen Beobachtungen Rückschlüsse auf das Modell ziehen?

*Beispiel 1.1 (Saskias Umfrage)* Versetzen wir uns in die Lage der Studentin Saskia, die für ihre Masterarbeit eine Befragung unter den 20.000 Studentinnen und Studenten ihrer Uni durchführen möchte. Von ihnen möchte Saskia wissen, ob sie mit ihrem Studium zufrieden sind. Dazu führt sie eine Online-Umfrage durch, bei der es nur eine Frage mit zwei Antwortmöglichkeiten gibt – man ist zufrieden oder nicht. Sie lässt den Studierenden die Umfrage über den internen Mailverteiler der Uni zukommen. Alle erhalten die Umfrage, aber nur ein Teil von ihnen macht sich die Mühe, die Umfrage zu beantworten. Welche Schlüsse kann Saskia aus ihrer Umfrage ziehen?

Nehmen wir an, die gewählten Fächer sind zu gleichen Teilen zufriedenstellend und nicht zufriedenstellend, das heißt 50% sind zufrieden, und die Umfrage wird wahrheitsgemäß beantwortet. Würde nur eine Studentin antworten, so entspräche das Ergebnis nicht der Wahrheit, denn es könnte nur eine einhundertprozentige (Un)Zufriedenheit ermittelt werden. Eine sogenannte *Stichprobengröße* von 1 ist

also zu niedrig. Je mehr Leute antworten, desto höher ist die *Wahrscheinlichkeit*, dass der Anteil der positiven Antworten in der Umfrage näher an den wahren 50% liegt. Nur wenn alle 20.000 Studierenden antworten würden, könnte Saskia sicher sein, den wahren Wert gefunden zu haben. Im Allgemeinen wird das Ergebnis von Saskias Umfrage also von Unsicherheit behaftet sein. Um die Aussagekraft des Ergebnisses einschätzen zu können, benötigen wir statistische Methoden. Andersherum wäre es für Saskia gut zu wissen, wie viele Antworten sie benötigt, um ein aussagekräftiges Ergebnis zu erzielen.

*Beispiel 1.2 (Wirksamkeit von Impfungen)* Ein Corona-Impfstoff wurde in einer ersten Studie bei  $n = 43\,500$  Teilnehmenden untersucht. Eine Hälfte bekam den Impfstoff verabreicht, die andere Hälfte als Kontrollgruppe ein Placebo. Innerhalb der darauffolgenden sechs Monate waren 4,03% in der Placebo-Gruppe an Corona erkrankt, aber nur 0,37% in der geimpften Gruppe. Dies wurde mit der Aussage „Das Erkrankungsrisiko wird durch Impfung um etwa 91% gesenkt“ veröffentlicht.

Dies sollten wir kritisch hinterfragen. Könnte es vielleicht nur Zufall gewesen sein, dass weniger Geimpfte erkrankten, ist der Unterschied also eventuell statistisch nicht signifikant? Mathematisch: wie wahrscheinlich ist eine solche Abweichung der Prozentzahlen bei gleicher Erkrankungswahrscheinlichkeit pro Patient in beiden Gruppen? Eine weitergehende Frage wäre, ob wichtige *Kovariablen* wie Geschlecht oder Alter die Erkrankungswahrscheinlichkeit in den beiden Gruppen beeinflussen, sodass sich vielleicht die Effizienz der Impfung detaillierter beschreiben lässt oder das Ergebnis sogar auf eine inhomogene Zusammensetzung der Gruppen zurückzuführen ist.

*Beispiel 1.3 (Happiness-Score)* Im *World Happiness Report* der Vereinten Nationen (UN) wird jährlich ein Glücksindex bzw. ein *Happiness-Score* zur Lebenszufriedenheit der Bevölkerung aus verschiedenen Ländern bestimmt. Wir wollen die Abhängigkeit des Happiness-Scores aus dem Jahr 2019 vom jeweiligen Pro-Kopf-Bruttoinlandsprodukt untersuchen und betrachten hierzu das *Modell*

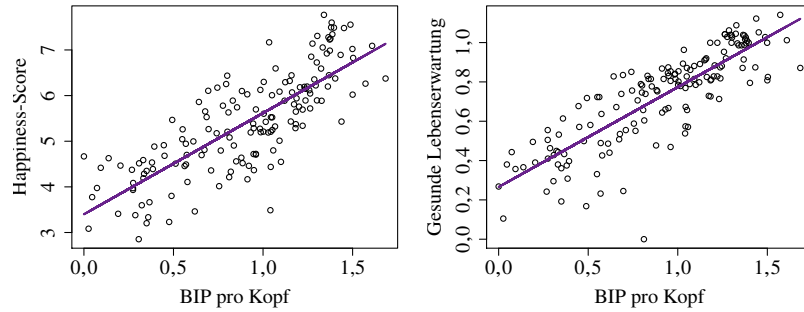
$$Y_i = aX_i + b + \varepsilon_i, \quad i = 1, \dots, 156, \quad (1.1)$$

wobei die zufälligen Störgrößen  $\varepsilon_i$  Unsicherheiten bei den Bevölkerungsumfragen sowie länderspezifische ökonomische/geographische/soziale Einflüsse etc. modellieren. Analog können wir versuchen, die Lebenserwartung bei guter Gesundheit durch das Bruttoinlandsprodukt zu erklären, siehe Abbildung 1.1. Plausible Annahmen an das Modell sind:

1.  $(\varepsilon_i)$  sind unabhängig (näherungsweise),
2.  $(\varepsilon_i)$  sind identisch verteilt,
3.  $\mathbb{E}[\varepsilon_i] = 0$  (kein systematischer Fehler),
4.  $\varepsilon_i$  normalverteilt (wegen des zentralen Grenzwertsatzes).

Naheliegende *Ziele* bzw. *Fragestellungen* sind:

1. Für welche Parameter  $a, b$  beschreibt das Modell (1.1) die gegebenen Daten am besten? Ein mögliches Schätzverfahren ist der *Kleinste-Quadrate-Schätzer*



**Abb. 1.1** Happiness-Score (*links*) und Lebenserwartung bei guter Gesundheit (*rechts*) in Abhängigkeit von Bruttoinlandsprodukt pro Kopf aus 156 Ländern auf Grundlage des World Happiness Reports von 2019. Die resultierenden Regressionsgeraden aus Beispiel 1.3 sind violett eingezeichnet.

$$(\hat{a}, \hat{b}) := \arg \min_{a, b \in \mathbb{R}} \sum_{i=1}^n (Y_i - aX_i - b)^2$$

(wir minimieren die Summe der quadrierten Residuen), den wir in Kapitel 2 kennenlernen werden. Mit  $\hat{a}$  und  $\hat{b}$  erhalten wir die *Regressionsgerade*

$$y = \hat{a}x + \hat{b}.$$

2. Sind die Modellannahmen erfüllt? Hierzu können Histogramme, Boxplots und Quantil-Quantil-Diagramme (QQ-Plots) der Residuen  $Y_i - \hat{a}X_i - \hat{b}$  verwendet werden.
3. Wenn wir die Verteilung von  $\hat{a}$  kennen (Verteilungsannahme an  $\varepsilon$  nötig!), können wir Intervalle der Form  $I = [\hat{a} - c, \hat{a} + c]$  für  $c > 0$  konstruieren, sodass der tatsächliche Parameter  $a$  mit vorgegebener Wahrscheinlichkeit in  $I$  liegt.
4. Wie kann man *testen*, ob das Bruttoinlandsprodukt tatsächlich einen Effekt auf die Lebenszufriedenheit hat, das heißt gilt die Hypothese  $H_0 : a = 0$  oder kann sie verworfen werden? Ein mögliches Verfahren ist, die Hypothese zu verwerfen, falls  $|\hat{a}| > c$  für einen kritischen Wert  $c > 0$ . Um einen sinnvollen Wert zu bestimmen, benötigen wir wieder Verteilungsannahmen an die Fehler ( $\varepsilon_i$ ).

Im Laufe der folgenden Kapitel werden solche und ähnliche Fragen beantwortet, doch vorher müssen wir statistische Modelle formal einführen.

**Definition 1.4** Sei  $\mathcal{X}$  die Menge aller möglichen Beobachtungen, der sogenannte **Stichprobenraum**. Ein messbarer Raum  $(\mathcal{X}, \mathcal{F})$ , versehen mit einer Familie  $(\mathbb{P}_\theta)_{\theta \in \Theta}$  von Wahrscheinlichkeitsmaßen mit einer beliebigen Parametermenge  $\Theta \neq \emptyset$ , heißt **statistisches Experiment** oder **statistisches Modell**.

Durch die Wahl der Familie von Wahrscheinlichkeitsmaßen im statistischen Modell wird zielgerichtet eine Veränderlichkeit oder Unsicherheit des Modells formalisiert. Je nachdem, welche statistische Aufgabe gegeben ist, wählt man statistische

Methoden, die wiederum Rückschlüsse auf das Modell erlauben. Zuletzt werden die im Experiment gewonnenen Daten in die Methoden eingesetzt, sodass man Informationen über den zugrunde liegenden Parameter  $\vartheta$  gewinnt, unter dem die Daten generiert wurden.

**Definition 1.5** Es sei  $(\Omega, \mathcal{A}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $(X, \mathcal{F})$  ein Messraum. Jede  $(\mathcal{A}, \mathcal{F})$ -messbare Abbildung  $X : \Omega \rightarrow X$  heißt **Beobachtung** oder **Statistik** mit Werten in  $(X, \mathcal{F})$  und induziert das statistische Modell  $(X, \mathcal{F}, (\mathbb{P}_\vartheta^X)_{\vartheta \in \Theta})$ . Sind die Beobachtungen  $X_1, \dots, X_n, n \in \mathbb{N}$ , für jedes  $\vartheta \in \Theta$  unabhängig und identisch verteilt (kurz i.i.d. für *independent and identically distributed*), so nennt man  $X_1, \dots, X_n$  eine **mathematische Stichprobe**.

In der Statistik konzentrieren wir uns auf die Beobachtungen und ihre Verteilung. Entsprechend der vorangegangenen Definition betrachten wir also das statistische Experiment  $(X, \mathcal{F}, (\mathbb{P}_\vartheta^X)_{\vartheta \in \Theta})$ , in dem die Werte von  $X : \Omega \rightarrow X$  liegen. Das statistische Modell  $(\Omega, \mathcal{A}, \mathbb{P}_{\vartheta \in \Theta})$ , in dem die Urbilder von  $X$  liegen, bleibt häufig unspezifiziert. Dies ist vollkommen analog zur Rolle der Zufallsvariablen in der Wahrscheinlichkeitstheorie. Der Einfachheit halber schreiben wir  $X \sim \mathbb{P}_\vartheta$  und nicht  $X \sim \mathbb{P}_\vartheta^X$ .

**Bemerkung 1.6** Für  $n \in \mathbb{N}$  sei  $X_1, \dots, X_n$  eine mathematische Stichprobe mit Werten in  $X$  und Randverteilung  $X_1 \sim \mathbb{P}_\vartheta$  mit Parameter  $\vartheta \in \Theta$ . Dann ist der Stichprobenvektor  $(X_1, \dots, X_n)$  gemäß dem Produktmaß  $\mathbb{P}_\vartheta^{\otimes n} = \bigotimes_{i=1}^n \mathbb{P}_\vartheta$  auf  $(X^n, \mathcal{F}^{\otimes n})$  verteilt, siehe Definition A.20.

**Beispiel 1.7 (Statistisches Modell, mathematische Stichprobe)** Kommen wir auf Beispiel 1.1 zurück. Die Studentin Saskia bekommt  $n$  Antworten auf ihre Umfrage zur Studienzufriedenheit unter  $N = 20.000$  Studentinnen und Studenten. Da wir davon ausgehen, dass eine bereits befragte Person nicht noch einmal an der Umfrage teilnimmt, bietet sich die hypergeometrische Verteilung  $\text{Hyp}(N, M, n)$  zur Modellierung an, wobei die Anzahl  $M$  der zufriedenen Studierenden der unbekannte Parameter ist. Falls deutlich weniger als  $N$  Antworten eingehen, können wir die hypergeometrische Verteilung durch die Binomialverteilung beziehungsweise mittels eines Vektors von Bernoulli-verteilten Zufallsvariablen approximieren (siehe Bemerkung A.48).

Folglich kann Saskia die eingegangenen Antworten  $X_1, \dots, X_n$  als unabhängige Bernoulli-verteilte Zufallsvariablen modellieren, wobei 0 für „unzufrieden“ und 1 für „zufrieden“ steht. Der Erfolgs- bzw. Zufriedenheitsparameter  $\vartheta = M/N \in [0, 1]$  ist gerade der Anteil der zufriedenen Studentinnen und Studenten. Wir erhalten das statistische Modell  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit

$$X = \{0, 1\}^n, \quad \mathcal{F} = \mathcal{P}(X), \quad \mathbb{P}_\vartheta = \text{Ber}(\vartheta)^{\otimes n} \quad \text{und} \quad \Theta = [0, 1].$$

**Beispiel 1.8 (Statistisches Modell II)** Wir modellieren nun die Situation aus Beispiel 1.2. Wir notieren den Versuchsausgang als einen 0-1-Vektor  $x$  der Länge  $n$ , wobei die ersten  $n/2$  Einträge die Ergebnisse in Gruppe 1 (Placebo) und die letzten  $n/2$  Einträge die Ergebnisse in Gruppe 2 (geimpft) kodieren mit '1'='erkrankt', '0'='nicht erkrankt'. Wir wählen also den Stichprobenraum  $X = \{0, 1\}^n$  mit der

$\sigma$ -Algebra  $\mathcal{F} = \mathcal{P}(\mathcal{X})$ . Wenn wir annehmen, dass alle Teilnehmenden in der Studie unabhängig voneinander erkranken oder nicht und die Erkrankungswahrscheinlichkeit  $p_1 \in [0, 1]$  in Gruppe 1 sowie  $p_2 \in [0, 1]$  in Gruppe 2 beträgt, so sind die  $n$  einzelnen Versuchsausgänge unabhängig Bernoulli-verteilt mit Parametern  $p_1$  bzw.  $p_2$ . Unter Verwendung des Produktmaßes und mit  $\Theta = [0, 1]^2$  ergibt sich daher das statistische Modell  $(\mathcal{X}, \mathcal{F}, (\text{Ber}(p_1)^{\otimes n/2} \otimes \text{Ber}(p_2)^{\otimes n/2})_{(p_1, p_2) \in \Theta})$ . Einfacher verständlich ist die Beschreibung, dass wir einen  $\mathcal{X}$ -wertigen Zufallsvektor beobachten, dessen Koordinaten  $X_1, \dots, X_n$  unabhängig sind mit  $X_i \sim \text{Ber}(p_1)$  für  $i = 1, \dots, n/2$  und  $X_i \sim \text{Ber}(p_2)$  für  $i = n/2 + 1, \dots, n$ , wobei  $(p_1, p_2) \in \Theta$  der unbekannte Parameter ist.

Typische Statistiken in dem Modell sind  $S_1 = \sum_{j=1}^{n/2} X_j \sim \text{Bin}(n/2, p_1)$  und  $S_2 = \sum_{j=n/2+1}^n X_j \sim \text{Bin}(n/2, p_2)$ . Es ist intuitiv, dass durch diese Zusammenfassung der Erkrankungszahlen in den einzelnen Gruppen keine statistische Information verloren geht.

Wie bereits diese einfachen Beispiele zeigen, bilden Modelle nicht die gesamte Wirklichkeit ab, sondern dienen immer auch der Vereinfachung. Folgendes Zitat sollte man sich daher immer vor Augen halten: „*Alle Modelle sind falsch, doch manche sind nützlich.*“ (George Box).

Viele statistische Fragestellungen kann man einem der drei Grundprobleme *Parameterschätzung*, *Hypothesentests* und *Konfidenzmengen* zuordnen. Diese werden im Folgenden eingeführt und später weiter vertieft. Weitere wichtige statistische Aufgaben wie Vorhersage oder Klassifikation werden sich durch leichte Modifikationen ergeben.

## 1.2 Parameterschätzung

Unser erstes Ziel ist, aufgrund von Beobachtungen den unbekannten Parameter  $\vartheta \in \Theta$  zu schätzen. Wir sagen „schätzen“ und nicht „bestimmen“, weil wir in den meisten Fällen mit Wahrscheinlichkeiten, Fehlern und unvollständigen Stichproben konfrontiert sind, wie in Beispiel 1.7 mit Saskia, in dem nur eine zufällige Teilmenge von  $n \leq 20.000$  Menschen Saskias Umfrage beantwortet. Wir können also nicht erwarten, den „wahren“ zugrunde liegenden Parameter, sollte es ihn überhaupt geben, zu finden. Stattdessen soll auf Grundlage der Daten ein Parameterwert gefunden werden, mit dem das resultierende Modell die Daten möglichst gut beschreibt.

### 1.2.1 Konstruktionsprinzipien

Ein Schätzer ist eine Funktion, manchmal auch *Schätzfunktion* genannt, die jeder realisierten Beobachtung  $x \in \mathcal{X}$  einen Parameter  $\vartheta \in \Theta$  zuordnet. Wenn wir auch die Schätzung von *abgeleiteten Parametern*  $\rho(\vartheta)$  zulassen, wobei  $\rho: \Theta \rightarrow \mathbb{R}^p$ , führt dies auf folgende Definition:

**Definition 1.9** Sei  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $\rho: \Theta \rightarrow \mathbb{R}^p$  ein (abgeleiteter)  $p$ -dimensionaler Parameter für  $p \in \mathbb{N}$ . Ein **Schätzer** ist eine messbare Abbildung  $\hat{\rho}: X \rightarrow \mathbb{R}^p$ . Gilt  $\mathbb{E}_\vartheta[\hat{\rho}] = \rho(\vartheta)$  für alle  $\vartheta \in \Theta$ , so heißt  $\hat{\rho}$  **unverzerrt** oder **erwartungstreu** (englisch: *unbiased*). Der Erwartungswert eines Zufallsvektors ist dabei koordinatenweise definiert.

Da für  $\Theta \subseteq \mathbb{R}^p$  obige Definition die Identität  $\rho(\vartheta) = \vartheta$  einschließt, ist die Betrachtung von abgeleiteten Parametern etwas allgemeiner als nur die Untersuchung von Schätzern des gesamten Parameters.

*Beispiel 1.10 (Schätzer, abgeleiteter Parameter – a)* Wir setzen Saskias Beispiel 1.7 fort und betrachten die mathematische Stichprobe  $X_1, \dots, X_n \sim \text{Ber}(\vartheta)$  mit Parameter  $\vartheta \in [0, 1]$ . Da der Parameterraum nur eindimensional ist, betrachten wir die Identität  $\rho(\vartheta) = \vartheta$ . Als Schätzer wählen wir den Anteil der zufriedenen Studierenden unter allen Antworten, das heißt

$$\hat{\rho}_n := \frac{1}{n} \sum_{i=1}^n X_i. \quad (1.2)$$

Würden wirklich alle 20.000 Studentinnen und Studenten an der Umfrage teilnehmen, würde  $\hat{\rho}_n$  exakt den Anteil der zufriedenen Studierenden ergeben. Unter unserer Modellannahme ist  $\hat{\rho}_n$  aber auch sonst sinnvoll: Es gilt

$$\mathbb{E}_\vartheta[\hat{\rho}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\vartheta[X_i] = \vartheta \quad \text{und} \quad \text{Var}_\vartheta(\hat{\rho}_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}_\vartheta(X_i) = \frac{\vartheta(1-\vartheta)}{n}.$$

Das heißt, der Schätzer  $\hat{\rho}_n$  ist erwartungstreu, und für größer werdenden Stichprobenumfang  $n$  sinkt seine Varianz.  $\hat{\rho}_n$  konzentriert sich also für wachsendes  $n$  um das wahre  $\vartheta$ . Beide Eigenschaften sind wünschenswert. Wir werden später Beispiele sehen, bei denen wir jedoch eine Abweichung des Erwartungswerts  $\mathbb{E}[\hat{\rho}_n]$  vom wahren Parameter, man spricht dann von einem *Bias*, in Kauf nehmen, um die Varianz zusätzlich zu reduzieren.

*Beispiel 1.11 (Schätzer, abgeleiteter Parameter – b)* Es sei  $X_1, \dots, X_n$  eine normalverteilte mathematische Stichprobe, das heißt  $X_1 \sim N(\mu, \sigma^2)$ . Der unbekannte Parameter ist

$$\vartheta = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty) =: \Theta.$$

Interessieren wir uns nur für den Erwartungswert  $\mu$ , dann müssen wir von dem zweidimensionalen Parameter  $\vartheta$  nur noch die erste Komponente schätzen. Hier kommt nun die den Parameter ableitende Funktion  $\rho$  ins Spiel, die  $\vartheta$  auf das reduziert, was uns interessiert. Wir definieren also  $\rho: \Theta \rightarrow \mathbb{R}, (\mu, \sigma^2) \mapsto \mu$ . Als Schätzer für  $\rho(\vartheta)$  verwenden wir wieder das arithmetische Mittel  $\bar{X}_n$ , auch *Stichprobenmittel* genannt, das hier ebenfalls erwartungstreu ist (dies zu zeigen sei den Leserinnen und Lesern überlassen)

$$\hat{\rho}_n := \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

Um die Güte eines Schätzers zu bestimmen, werden *Verlustfunktionen* und deren zugehöriges *Risiko* verwendet. Diese messen den (erwarteten) Abstand zwischen geschätztem und wahren Parameter.

**Definition 1.12** Eine Funktion  $\ell: \Theta \times \mathbb{R}^p \rightarrow \mathbb{R}_+$  heißt **Verlustfunktion**, falls  $\ell(\vartheta, \cdot)$  für jedes  $\vartheta \in \Theta$  messbar ist. Der erwartete Verlust  $R(\vartheta, \hat{\rho}) := \mathbb{E}_\vartheta[\ell(\vartheta, \hat{\rho})]$  eines Schätzers  $\hat{\rho}$  heißt **Risiko**. Typische Verlustfunktionen sind

1. der **0-1-Verlust**  $\ell(\vartheta, r) = \mathbb{1}_{\{r \neq \rho(\vartheta)\}}$ ,
2. der **absolute Verlust**  $\ell(\vartheta, r) = |r - \rho(\vartheta)|$  (euklidischer Abstand im  $\mathbb{R}^p$ ) sowie
3. der **quadratische Verlust**  $\ell(\vartheta, r) = |r - \rho(\vartheta)|^2$ .

Das Risiko bezüglich des quadratischen Verlusts nennt man auch *quadratisches Risiko* bzw. *mittleren quadratischen Fehler* (englisch: *mean squared error*, kurz: MSE). Dieser lässt sich in Bias und Varianz zerlegen. Wir erinnern hierfür daran, dass der Erwartungswert eines Zufallsvektors  $X = (X_1, \dots, X_d)^\top \in \mathbb{R}^d$  koefizientenweise gebildet wird, das heißt  $\mathbb{E}[X] := (\mathbb{E}[X_1], \dots, \mathbb{E}[X_d])^\top$ , und die Varianz durch

$$\text{Var}(X) := \mathbb{E}[|X - \mathbb{E}[X]|^2] = \sum_{i=1}^d \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] = \sum_{i=1}^d \text{Var}(X_i)$$

definiert wird.

**Lemma 1.13 (Bias-Varianz-Zerlegung, statistischer Pythagoras)** *Es sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $\hat{\rho}: \mathcal{X} \rightarrow \mathbb{R}^p$  ein Schätzer des Parameters  $\rho(\vartheta)$  mit  $\mathbb{E}_\vartheta[|\hat{\rho}|^2] < \infty$  für alle  $\vartheta \in \Theta$ . Dann gilt für das quadratische Risiko*

$$\mathbb{E}_\vartheta[|\hat{\rho} - \rho(\vartheta)|^2] = \text{Var}_\vartheta(\hat{\rho}) + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2 \quad \text{für alle } \vartheta \in \Theta.$$

Die mittlere Abweichung des Schätzers vom wahren Wert  $\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)$  wird **Bias** genannt.

**Beweis** Für Vektoren  $a, b \in \mathbb{R}^p$  gilt  $|a+b|^2 = |a|^2 + 2\langle a, b \rangle + |b|^2$ . Aus der Linearität des Erwartungswerts folgt für alle  $\vartheta \in \Theta$ :

$$\begin{aligned} \mathbb{E}_\vartheta[|\hat{\rho} - \rho(\vartheta)|^2] &= \mathbb{E}_\vartheta[|\hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}] + \mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2] \\ &= \mathbb{E}_\vartheta[|\hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}]|^2] + 2\mathbb{E}_\vartheta[\langle \hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}], \mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta) \rangle] \\ &\quad + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2 \\ &= \text{Var}_\vartheta(\hat{\rho}) + 2\langle \mathbb{E}_\vartheta[\hat{\rho}] - \mathbb{E}_\vartheta[\hat{\rho}], \mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta) \rangle \\ &\quad + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2 \\ &= \text{Var}_\vartheta(\hat{\rho}) + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2 \end{aligned}$$

Damit ist die Behauptung gezeigt. □

*Beispiel 1.14 (Quadratisches Risiko)* Im Saskia-Beispiel 1.10 betrachten wir nun den Schätzer  $\tilde{\rho}_n := (\sum_{i=1}^n X_i + 1)/(n + 2)$ . Dieser hat den Bias

$$\mathbb{E}_\vartheta[\tilde{\rho}_n] - \vartheta = \frac{1 - 2\vartheta}{n + 2}$$

und die Varianz

$$\text{Var}_\vartheta(\tilde{\rho}_n) = \frac{n\vartheta(1 - \vartheta)}{(n + 2)^2}.$$

Damit gilt  $\mathbb{E}_\vartheta[\tilde{\rho}_n] = \rho(\vartheta)$  ausschließlich für  $\vartheta = 1/2$ , sodass  $\tilde{\rho}_n$  kein unverzerrter Schätzer ist. Er besitzt aber einen kleineren quadratischen Fehler als  $\hat{\rho}_n$ , wenn  $|\vartheta - 1/2| \leq 1/\sqrt{8}$ .

Obwohl wir nur wenig Asymptotik behandeln, also das Verhalten der Schätzer bei Stichprobenumfängen  $n \rightarrow \infty$ , seien noch zwei weitere wichtige Grundbegriffe erwähnt, die ebenfalls als Gütekriterien eines Schätzers dienen.

**Definition 1.15** Für jedes  $n \in \mathbb{N}$  sei  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \mathbb{P}_\vartheta$  eine mathematische Stichprobe. Dann heißt eine Folge von Schätzern  $\hat{\rho}_n = \hat{\rho}_n(X_1, \dots, X_n)$  des abgeleiteten Parameters  $\rho(\vartheta)$  **konsistent**, falls für alle  $\vartheta \in \Theta$

$$\hat{\rho}_n \xrightarrow{\mathbb{P}_\vartheta} \rho(\vartheta) \quad \text{für } n \rightarrow \infty.$$

Falls  $\mathbb{E}_\vartheta[|\hat{\rho}_n|^2] < \infty$  für alle  $\vartheta \in \Theta$  und unter jedem  $\mathbb{P}_\vartheta$

$$r_n^{-1}(\hat{\rho}_n - \rho(\vartheta)) \xrightarrow{d} N(\mu_\vartheta, \sigma_\vartheta^2) \quad \text{für } n \rightarrow \infty$$

gilt, nennen wir  $\hat{\rho}_n$  **asymptotisch normalverteilt** mit Konvergenzrate  $r_n \rightarrow \infty$ , asymptotischem Bias  $\mu_\vartheta \in \mathbb{R}$  und asymptotischer Varianz  $\sigma_\vartheta^2$ .

Warum ist es sinnvoll, einen Schätzer (genauer eine Folge von Schätzern) in diesem Fall „konsistent“ zu nennen? „Konsistent“ bedeutet in der Philosophie, genauer in der Logik, „zusammenhängend in der Gedankenführung“ und tatsächlich stellt obige Definition einen (asymptotischen) Zusammenhang zwischen dem Schätzer  $\hat{\rho}_n$  und dem zu schätzenden Parameter  $\rho(\vartheta)$  her.

Ein stärkeres Kriterium beschreibt die asymptotische Verteilung eines Schätzers. Aufgrund des zentralen Grenzwertsatzes sind viele Schätzer asymptotisch normalverteilt, so auch in Beispiel 1.10. Daher kommt der Untersuchung von statistischen Modellen unter Normalverteilungsannahme eine besondere Bedeutung zu.

Bisher haben wir unsere Beispielschätzer eher intuitiv gefunden. Stattdessen wollen wir nun zwei wichtige Konstruktionsprinzipien einführen, die *Momentenmethode* und die *Maximum-Likelihood-Schätzung*, die in einem sehr allgemeinen Rahmen eingesetzt werden können. Beide Ansätze liefern häufig gute Schätzmethoden.

Wir beginnen mit der *Momentenmethode*. Wie der Name suggeriert, wollen wir die Momente von Zufallsvariablen verwenden, um den unbekannten Parameter zu schätzen. Betrachten wir ein statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit  $\Theta \subseteq \mathbb{R}^p$ , so

hängen die Momente  $m_k(\vartheta) := \mathbb{E}_\vartheta[X^k]$  einer Beobachtung  $X \sim \mathbb{P}_\vartheta$  im Allgemeinen vom Parameter  $\vartheta$  ab. Die Idee ist nun, die  $p$  Komponenten von  $\vartheta$  zu rekonstruieren, indem man ein Gleichungssystem mit  $p$  Gleichungen basierend auf den (ersten)  $p$  Momenten  $m_k(\vartheta)$  löst. Für den Fall, dass dieses Gleichungssystem eine eindeutige Lösung besitzt, muss man nur noch die Momente schätzen. Haben wir  $n$  i.i.d. Beobachtungen  $X_1, \dots, X_n \sim \mathbb{P}_\vartheta$  zur Verfügung, konvergiert nach dem Gesetz der großen Zahlen das  $k$ -te Stichprobenmoment  $\widehat{m}_k := \frac{1}{n} \sum_{j=1}^n X_j^k$  für  $n \rightarrow \infty$  gegen  $m_k(\vartheta)$ .

**Methode 1.16 (Momentenmethode, Momentenschätzer)** Sei  $\vartheta \in \mathbb{R}^p$  der unbekannte Parameter einer mathematischen Stichprobe  $X_1, \dots, X_n$  reeller Zufallsvariablen. Für  $k \in \mathbb{N}$  bezeichne

$$m_k := m_k(\vartheta) := \mathbb{E}_\vartheta[X_1^k] \quad \text{bzw.} \quad \widehat{m}_k := \frac{1}{n} \sum_{j=1}^n X_j^k$$

das  $k$ -te Moment von  $X_1$ , sofern es existiert ( $\mathbb{E}_\vartheta[|X_1|^k] < \infty$ ), bzw. das  $k$ -te Stichprobenmoment. Für gewählte Indizes  $1 \leq k_1 < k_2 < \dots < k_p$  mit  $\mathbb{E}_\vartheta[|X_1|^{k_p}] < \infty$  ergibt sich der (klassische) **Momentenschätzer**  $\widehat{\vartheta}$  von  $\vartheta$  als die Lösung der  $p$  Gleichungen

$$m_{k_1}(\widehat{\vartheta}) = \widehat{m}_{k_1}, \quad m_{k_2}(\widehat{\vartheta}) = \widehat{m}_{k_2}, \quad \dots, \quad m_{k_p}(\widehat{\vartheta}) = \widehat{m}_{k_p}. \quad (1.3)$$

In vielen Fällen ist das Gleichungssystem mit  $p$  Gleichungen eindeutig lösbar, aber es gibt Beispiele, in denen keine eindeutig bestimmte Lösung existiert. Dies und die resultierenden Eigenschaften, beispielsweise das quadratische Risiko, bestimmen die Wahl der Exponenten  $k_1, \dots, k_p$ . Ganz allgemein stellt sich die Frage, wann die Folge der Momente  $(m_k)_{k \geq 1}$ , falls sie existiert, die Verteilung  $\mathbb{P}_\vartheta$  eindeutig bestimmt. Eine mögliche Antwort auf dieses sogenannte *Momentenproblem* liefert der Satz A.30. Über die Monome  $X_1^k$  hinausgehend, lässt sich die Momentenmethode auch direkt auf Erwartungswerte allgemeinerer Funktionale  $f(X_1)$  oder auf zentrierte Momente wie die Varianz  $\text{Var}_\vartheta(X_1)$  mit empirischer Varianz  $\widehat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \widehat{m}_1)^2$  verallgemeinern.

**Beispiel 1.17 (Momentenschätzer)** Sei  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$ . Dann ist  $m_1 = \mathbb{E}_{\mu, \sigma^2}[X_1] = \mu$  und  $m_2 = \mathbb{E}_{\mu, \sigma^2}[X_1^2] = \text{Var}_{\mu, \sigma^2}(X_1) + \mathbb{E}_{\mu, \sigma^2}[X_1]^2 = \sigma^2 + \mu^2$ . Wir erhalten das Gleichungssystem

$$\widehat{\mu} = \frac{1}{n} \sum_{j=1}^n X_j \quad \text{und} \quad \widehat{\sigma}^2 + \widehat{\mu}^2 = \frac{1}{n} \sum_{j=1}^n X_j^2.$$

Mithilfe des Stichprobenmittels  $\overline{X}_n = \frac{1}{n} \sum_{j=1}^n X_j$  erhalten wir die Lösung

$$\widehat{\mu} = \overline{X}_n, \quad \widehat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n X_j^2 - \left( \frac{1}{n} \sum_{j=1}^n X_j \right)^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \overline{X}_n)^2.$$

**Kurzbiografie (Karl und Egon Pearson)** Karl Pearson wurde 1857 in London geboren. Er studierte in Cambridge, Heidelberg und Berlin Mathematik, Physik, Deutsche Literatur, Biologie, Philosophie, Jura und Geschichte. Den Großteil seiner Laufbahn verbrachte er am Londoner University College. Unter seinen Arbeiten in zahlreichen wissenschaftlichen Disziplinen sind die „Mathematical Contributions to the Theory of Evolution“ sein wertvollster Beitrag zur Statistik. Unter anderem entwickelte er darin die *Momentenmethode*, den Korrelationskoeffizienten (nach einer ähnlichen Idee von Francis Galton) und den  $\chi^2$ -Test der statistischen Signifikanz. Karl Pearson starb 1936 in Coldharbour, Surrey.

Auch Karl Pearsons Sohn Egon Sharpe Pearson (geboren 1895 in Hampstead, gestorben 1980 in Sussex) wurde einer der führenden britischen Statistiker. Egon Pearson folgte seinem Vater an das University College London, und als Karl Pearson in den Ruhestand ging, wurde sein Lehrstuhl auf seinen Sohn Egon sowie Ronald Aylmer Fisher aufgeteilt. Egon Pearson ist insbesondere für das *Neyman-Pearson-Lemma* bekannt.

Für die zweite Konstruktionsmethode von Schätzern benötigen wir etwas mehr Struktur, die wir auch im weiteren Verlauf immer wieder aufgreifen. Wir fordern Absolutstetigkeit der Wahrscheinlichkeitsmaße (siehe Definition A.32), sodass uns der Satz von Radon-Nikodym (siehe Satz A.33) Dichten liefert.

**Definition 1.18** Ein statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  heißt **dominiert**, falls es ein  $\sigma$ -endliches Maß  $\mu$  gibt, sodass  $\mathbb{P}_\vartheta$  absolutstetig bezüglich  $\mu$  ist für alle  $\vartheta \in \Theta$ . Für jedes  $\vartheta \in \Theta$  bezeichnen wir die Radon-Nikodym-Dichte von  $\mathbb{P}_\vartheta$  bezüglich  $\mu$  mit

$$L(\vartheta, x) := \frac{d\mathbb{P}_\vartheta}{d\mu}(x), \quad \vartheta \in \Theta, x \in \mathcal{X},$$

und nennen sie **Likelihood-Funktion** oder kurz **Likelihood**. Äquivalent gilt  $L(\vartheta, x) = p_\vartheta(x)$  immer, wenn die  $\mu$ -Dichte  $p_\vartheta$  von  $\mathbb{P}_\vartheta$  für alle  $\vartheta \in \Theta$  existiert.

Damit ist  $x \mapsto L(\vartheta, x)$  eindeutig bis auf eine  $\mu$ -Nullmenge, die dann auch eine  $\mathbb{P}_{\vartheta'}$ -Nullmenge für jedes  $\vartheta' \in \Theta$  ist. Da wir aufgrund von gegebenen Daten  $X$  den zugrunde liegenden Parameter  $\vartheta$  schätzen wollen, wird in der Statistik die Likelihood-Funktion als durch  $x$  parametrisierte Funktion in  $\vartheta$  aufgefasst. Somit ist  $\vartheta \mapsto L(\vartheta) = L(\vartheta, X)$  eine zufällige Funktion, der sogenannte Likelihood-Prozess.

*Beispiel 1.19 (Likelihood-Funktionen)*

1. Betrachte das statistische Modell  $(\mathbb{R}, \mathcal{B}(\mathbb{R}), (\mathbb{N}(\mu, \sigma^2))_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)})$ . Als dominierendes Maß kann das Lebesgue-Maß gewählt werden, sodass die Likelihood-Funktion durch die Lebesgue-Dichte der Normalverteilung gegeben ist:  $L((\mu, \sigma^2), x) = (2\pi\sigma^2)^{-1/2} e^{-(x-\mu)^2/(2\sigma^2)}$ ,  $x \in \mathbb{R}$ .
2. Jedes statistische Modell auf dem Stichprobenraum  $(\mathbb{N}, \mathcal{P}(\mathbb{N}))$  oder allgemeiner auf einem abzählbaren Raum  $(\mathcal{X}, \mathcal{P}(\mathcal{X}))$  ist vom Zählmaß dominiert. Damit ist die Likelihood-Funktion durch die Zähldichte gegeben.
3. Ist  $\Theta = \{\vartheta_1, \vartheta_2, \dots\}$  abzählbar, so ist  $\mu = \sum_i c_i \mathbb{P}_{\vartheta_i}$  mit  $c_i > 0$  und  $\sum_i c_i = 1$  ein dominierendes Maß.

**Methode 1.20 (Maximum-Likelihood-Schätzer)** Für ein dominiertes statistisches Modell  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit Likelihood-Funktion  $L(\vartheta, x)$  heißt eine Statistik  $\hat{\vartheta}: X \rightarrow \Theta$  ( $\Theta$  trage eine  $\sigma$ -Algebra) **Maximum-Likelihood-Schätzer** (englisch: *maximum likelihood estimator*, kurz: MLE), falls gilt:

$$L(\hat{\vartheta}(x), x) = \sup_{\vartheta \in \Theta} L(\vartheta, x) \quad \text{für } \mu\text{-f.a. } x \in X$$

Die Grundidee der Maximum-Likelihood-Schätzung ist im Fall diskreter Beobachtungen intuitiv: Wir wählen denjenigen Parameter  $\hat{\vartheta}(x)$  aus, unter dem die Beobachtung  $X = x$  die größte Wahrscheinlichkeit besitzt. Als einfaches Beispiel betrachte die Beobachtung  $X \sim \text{Bin}(2, \vartheta)$  mit  $\vartheta \in [0, 1]$  unbekannt. Mit dem Zählmaß als dominierendem Maß gilt dann  $L(\vartheta, 0) = (1 - \vartheta)^2$ ,  $L(\vartheta, 1) = 2\vartheta(1 - \vartheta)$ ,  $L(\vartheta, 2) = \vartheta^2$ . Der MLE ist also eindeutig gegeben durch  $\hat{\vartheta}(0) = 0$ ,  $\hat{\vartheta}(1) = 1/2$ ,  $\hat{\vartheta}(2) = 1$ , das heißt  $\hat{\vartheta} = X/2$ .

Im allgemeinen Fall, zum Beispiel bei einer Likelihood-Funktion bezüglich des Lebesgue-Maßes, führt das Maximum-Likelihood-Prinzip sehr häufig, aber durchaus nicht immer, zu vernünftigen Schätzern. Es sei der Leserin ans Herz gelegt, nachzuweisen, dass die MLE-Eigenschaft nicht vom dominierenden Maß  $\mu$  abhängt (Aufgabe 1.3).

**Kurzbiografie (Ronald Aylmer Fisher)** Sir Ronald Aylmer Fisher wurde 1890 in London geboren. Von 1909 bis 1912 studierte er in Cambridge. Fisher lehrte zunächst an verschiedenen öffentlichen Schulen, bevor er anfang an der Rothamsted Experimental Station zu arbeiten, wo er in den Feldern Genetik, Evolution und Statistik forschte. Nach Pearsons Ausscheiden übernahm er seine Professur am University College London. 1943 besetzte er dann den Lehrstuhl für Genetik in Cambridge. Fisher prägte die theoretische Basis der Statistik und trug maßgeblich zu ihrem heutigen Wesen bei. Unter anderem führte er die Likelihood-Funktionen und den *Maximum-Likelihood-Schätzer*, die *analysis of variance* (ANOVA), das Konzept der Suffizienz, die nach ihm benannte Fisher-Information und vieles mehr ein. Ronald Aylmer Fisher starb 1962 in Adelaide, Australien.

**Beispiel 1.21 (Maximum-Likelihood-Schätzer – a)** Betrachten wir wieder eine mathematische Stichprobe  $X_1, \dots, X_n$  normalverteilter Zufallsvariablen. Dann ist das statistische Modell  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), (\mathbb{P}_{\mu, \sigma^2}^{\otimes n})_{(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)})$  mit  $\mathbb{P}_{\mu, \sigma^2} = N(\mu, \sigma^2)$  vom Lebesgue-Maß auf  $\mathbb{R}^n$  dominiert mit Likelihood-Funktion

$$L(\mu, \sigma^2; x) = (2\pi\sigma^2)^{-n/2} \prod_{j=1}^n \exp\left(-\frac{(x_j - \mu)^2}{2\sigma^2}\right), \quad x = (x_1, \dots, x_n) \in \mathbb{R}^n.$$

Um den Maximum-Likelihood-Schätzer zu berechnen, nutzen wir, dass Extremstellen unter monotonen Transformationen erhalten bleiben. Die Anwendung des Logarithmus auf die Likelihood-Funktion ist eine solche monotone Transformation, die den Vorteil hat, dass aus dem Produkt eine Summe wird (schon Fisher verwendete diesen Trick). Diese sogenannte *Loglikelihood-Funktion* ist hier

$$l(\mu, \sigma^2; x) := \log L(\mu, \sigma^2; x) = -\frac{n}{2}(\log(2\pi) + \log \sigma^2) - \sum_{j=1}^n \frac{(x_j - \mu)^2}{2\sigma^2}.$$

An einer Maximalstelle von  $\ell$  verschwinden die ersten Ableitungen und wir erhalten

$$0 = \sigma^{-2} \sum_{j=1}^n (x_j - \mu), \quad \frac{n}{2\sigma^2} = \frac{1}{2\sigma^4} \sum_{j=1}^n (x_j - \mu)^2.$$

Umstellen der ersten Gleichung nach  $\mu$  liefert  $\hat{\mu} = \bar{X}_n$  und Einsetzen in die zweite Gleichung ergibt  $\hat{\sigma}^2 = n^{-1} \sum_j (X_j - \bar{X}_n)^2$ . Es ist leicht nachzuprüfen, dass  $\hat{\mu}$  und  $\hat{\sigma}^2$  tatsächlich das Maximierungsproblem lösen (und messbar sind). In diesem Fall stimmt der Maximum-Likelihood-Schätzer also mit dem Momentenschätzer aus Beispiel 1.17 überein.

*Beispiel 1.22 (Maximum-Likelihood-Schätzer – b)* Es sei  $X_1, \dots, X_n \sim \text{Pois}(\lambda)$  eine Poisson-verteilte mathematische Stichprobe mit Parameter  $\lambda > 0$ , das heißt  $\mathbb{P}_\lambda(X_1 = k) = \lambda^k e^{-\lambda} / k!$ . Dann ist die gemeinsame Verteilung gegeben durch

$$\mathbb{P}_\lambda(X_1 = k_1, \dots, X_n = k_n) = \frac{\lambda^{k_1 + \dots + k_n} e^{-n\lambda}}{k_1! \dots k_n!}, \quad k_1, \dots, k_n \in \mathbb{N}_0,$$

und wird vom Zählmaß dominiert. Ableiten der Loglikelihood-Funktion nach  $\lambda$  und Nullsetzen führt auf den Maximum-Likelihood-Schätzer  $\hat{\lambda} = \bar{X}_n$ . Der zugehörige Nachweis einer hinreichenden Bedingung sei der Leserin überlassen.

*Beispiel 1.23 (Maximum-Likelihood-Schätzer – c)* Erinnern wir uns an das Saskia-Beispiel 1.7. Im fortführenden Beispiel 1.10 hatten wir im statistischen Modell  $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\text{Ber}(\vartheta)^{\otimes n})_{\vartheta \in [0, 1]})$  den Schätzer  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  gewählt. Zu welchem Schätzer führt der Maximum-Likelihood-Ansatz? Die Likelihood-Funktion ist hier gegeben durch

$$L(\vartheta, x) = \prod_{i=1}^n \vartheta^{x_i} (1 - \vartheta)^{1-x_i}, \quad x = (x_1, \dots, x_n) \in \{0, 1\}^n, \quad \vartheta \in [0, 1],$$

sodass wir für  $\vartheta \in (0, 1)$  die Loglikelihood-Funktion

$$l(\vartheta, x) = \left( \sum_{i=1}^n x_i \right) \log(\vartheta) + \left( n - \sum_{i=1}^n x_i \right) \log(1 - \vartheta)$$

erhalten. Dabei fällt auf, dass die Loglikelihood-Funktion nur von der Statistik  $s = \sum_{i=1}^n x_i$  abhängt. Ableiten von  $l(\cdot, x)$  ergibt  $\frac{\partial}{\partial \vartheta} l(\vartheta, x) = \frac{s}{\vartheta} + \frac{s-n}{1-\vartheta}$  und dessen Nullstelle ist gegeben durch  $\hat{\vartheta} = \frac{s}{n} = \frac{1}{n} \sum_{i=1}^n x_i$ . Eine weitere Ableitung von  $l(\cdot, x)$  liefert die hinreichenden Bedingung und die Fälle  $\vartheta \in \{0, 1\}$  diskutiert man leicht separat. Insgesamt ist  $\hat{\vartheta} = \frac{1}{n} \sum_{i=1}^n X_i$  tatsächlich der Maximum-Likelihood-Schätzer.

### 1.2.2 Minimax- und Bayes-Ansatz

Wir haben nun verschiedene Schätzmethoden, wie den Maximum-Likelihood-Schätzer oder die Momentenmethode kennengelernt. Natürlich gibt es weitere Konstruktionen und zahlreiche Variationen. Welche Konstruktionsmethode sollte anhand des gegebenen Schätzproblems ausgewählt werden?

Betrachten wir ein statistisches Modell  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit abgeleitetem Parameter  $\rho: \Theta \rightarrow \mathbb{R}^p$ . Weiter sei eine Verlustfunktion  $\ell$  gegeben. Als mögliches Vergleichskriterium hatten wir bereits die Risikofunktion  $R(\vartheta, \hat{\rho}) = \mathbb{E}_\vartheta[\ell(\vartheta, \hat{\rho})]$  eines Schätzers  $\hat{\rho}$  eingeführt. Man beachte jedoch folgendes Beispiel:

*Beispiel 1.24 (Risiken von Schätzern – a)* Es sei  $X \sim N(\mu, 1)$ ,  $\mu \in \mathbb{R}$ , und  $\ell(\mu, \hat{\mu}) := (\hat{\mu} - \mu)^2$ . Wir betrachten die beiden Schätzer  $\hat{\mu}_1 := X$  und  $\hat{\mu}_2 := 5$ . Die Risiken sind dann gegeben durch

$$R(\mu, \hat{\mu}_1) = \mathbb{E}_\vartheta[(X - \mu)^2] = 1 \quad \text{und} \quad R(\mu, \hat{\mu}_2) = (5 - \mu)^2.$$

Damit hat  $\hat{\mu}_1$  genau dann ein kleineres Risiko als  $\hat{\mu}_2$ , wenn  $\mu \notin [4, 6]$ . Welchen Schätzer sollen wir nun wählen in Anbetracht dessen, dass  $\mu$  unbekannt ist? Ein möglicher Ansatz ist, die maximalen Risiken der Schätzer über alle  $\mu \in \Theta$  zu bestimmen und dann den Schätzer zu wählen, der das kleinste maximale Risiko besitzt. Da  $R(\mu, \hat{\mu}_1) = 1$  konstant in  $\mu \in \mathbb{R}$  ist, während  $\mu \mapsto R(\mu, \hat{\mu}_2)$  eine nach oben geöffnete Parabel beschreibt und damit beliebig groß werden kann, würden wir den Schätzer  $\hat{\mu}_1$  verwenden.

**Definition 1.25** Im statistischen Modell  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit abgeleitetem Parameter  $\rho: \Theta \rightarrow \mathbb{R}^p$  und Verlustfunktion  $\ell$ , heißt ein Schätzer  $\hat{\rho}$  **minimax**, falls für das zugehörige Risiko  $R$

$$\sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho}) = \inf_{\tilde{\rho}} \sup_{\vartheta \in \Theta} R(\vartheta, \tilde{\rho})$$

gilt, wobei sich das Infimum über alle Schätzer (das heißt messbaren Funktionen)  $\tilde{\rho}: X \rightarrow \mathbb{R}^p$  erstreckt.

Anstatt den maximal zu erwartenden Verlust zu minimieren, könnten wir die Risiken abhängig von  $\vartheta$  gewichten und diese vergleichen. Der Schätzer mit dem geringsten gemittelten Risiko wird dann gewählt. Die nötige Struktur dafür gibt uns folgende Definition.

**Definition 1.26** Der Parameterraum  $\Theta$  trage eine  $\sigma$ -Algebra  $\mathcal{F}_\Theta$ ,  $\vartheta \mapsto \mathbb{P}_\vartheta(B)$  sei messbar für alle  $B \in \mathcal{F}$ , und die Verlustfunktion  $\ell$  sei produktmessbar, das heißt

$$\ell: (\Theta \times \mathbb{R}^p, \mathcal{F}_\Theta \otimes \mathcal{B}(\mathbb{R}^p)) \rightarrow (\mathbb{R}_+, \mathcal{B}(\mathbb{R}_+)) \text{ messbar.}$$

Als **a-priori-Verteilung**  $\pi$  des Parameters  $\vartheta$  bezeichnen wir ein Wahrscheinlichkeitsmaß auf  $(\Theta, \mathcal{F}_\Theta)$ . Das zu  $\pi$  assoziierte **Bayes-Risiko** eines Schätzers  $\hat{\rho}$  ist

$$R_\pi(\hat{\rho}) := \mathbb{E}_\pi[R(\vartheta, \hat{\rho})] = \int_\Theta \int_X \ell(\vartheta, \hat{\rho}(x)) \mathbb{P}_\vartheta(dx) \pi(d\vartheta).$$

Der Schätzer  $\hat{\rho}$  heißt **Bayes-Schätzer** oder **Bayes-optimal** (bezüglich  $\pi$ ), falls

$$R_{\pi}(\hat{\rho}) = \inf_{\tilde{\rho}} R_{\pi}(\tilde{\rho}),$$

wobei sich das Infimum über alle Schätzer (das heißt messbaren Funktionen)  $\tilde{\rho}: \mathcal{X} \rightarrow \mathbb{R}^P$  erstreckt.

A-priori-Verteilung oder englisch *prior distribution* ist insofern ein passender Name, als dass diese Verteilung gewählt wird, bevor die Daten beobachtet werden.

**Bemerkung 1.27 (Bayes-Risiko)** Das Bayes-Risiko kann auch als insgesamt zu erwartender Verlust in folgendem Sinne verstanden werden: Definiere  $\Omega := \mathcal{X} \times \Theta$  und die gemeinsame Verteilung von Beobachtung und Parameter  $\tilde{\mathbb{P}}$  auf  $(\mathcal{X} \times \Theta, \mathcal{F} \otimes \mathcal{F}_{\Theta})$  gemäß  $\tilde{\mathbb{P}}(dx, d\theta) = \mathbb{P}_{\theta}(dx)\pi(d\theta)$ . Bezeichnen  $X$  und  $T$  die Koordinatenprojektionen von  $\Omega$  auf  $\mathcal{X}$  bzw.  $\Theta$ , dann gilt  $R_{\pi}(\hat{\rho}) = \mathbb{E}_{\tilde{\mathbb{P}}}[\ell(T, \hat{\rho}(X))]$ .

**Beispiel 1.28 (Risiken von Schätzern – b)** Wir beobachten  $X \sim \text{Bin}(n, \vartheta)$  mit  $\vartheta \in [0, 1]$  unbekannt und betrachten die Schätzer  $\hat{\vartheta}_1 = X/n$ ,  $\hat{\vartheta}_2 = 0,5$  (unabhängig von den Daten!). Dann gilt für den mittleren quadratischen Fehler  $R(\vartheta, \hat{\vartheta}_1) = \vartheta(1-\vartheta)/n$  und  $R(\vartheta, \hat{\vartheta}_2) = (\vartheta - 0,5)^2$ . Das maximale Risiko beträgt  $1/(4n)$  für  $\hat{\vartheta}_1$  und  $1/4$  für  $\hat{\vartheta}_2$ , sodass wir unter diesem Kriterium  $\hat{\vartheta}_1$  für  $n \geq 2$  bevorzugen. Unter der a-priori-Verteilung  $\pi = U([0, 1])$  ist  $\vartheta$  gleichmäßig auf  $[0, 1]$  verteilt, und Integration ergibt die Bayesrisiken  $R_{\pi}(\hat{\vartheta}_1) = 1/(6n)$  und  $R_{\pi}(\hat{\vartheta}_2) = 1/12$ . Unter diesem Bayesrisiko ziehen wir also nur für  $n \geq 3$  den Schätzer  $\hat{\vartheta}_1$  vor.

Die a-priori-Verteilung  $\pi$  wird manchmal auch als subjektive Einschätzung der Verteilung des zugrunde liegenden Parameters interpretiert. Nachdem eine Beobachtung gemacht wurde, möchte man die Parameterverteilung mithilfe der zusätzlichen Information aktualisieren („Bayesian updating“). Um diesen Vorgang formal zu beschreiben, erinnern wir uns an die bedingten Dichten aus Definition A.34 und die Bayes-Formel (Satz A.35).

**Definition 1.29** Sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\theta})_{\theta \in \Theta})$  ein von  $\mu$  dominiertes statistisches Modell mit Dichten  $f^{X|T=\theta} := \frac{d\mathbb{P}_{\theta}}{d\mu}$ . Sei  $\pi$  eine a-priori-Verteilung auf  $(\Theta, \mathcal{F}_{\Theta})$  mit Dichte  $f^T$  bezüglich eines Maßes  $\nu$ . Ist  $f^{X|T=\cdot}: \mathcal{X} \times \Theta \rightarrow \mathbb{R}_+$  ( $\mathcal{F} \otimes \mathcal{F}_{\Theta}$ )-messbar, dann ist die **a-posteriori-Verteilung** des Parameters, gegeben die Beobachtung  $X = x$ , definiert durch die  $\nu$ -Dichte

$$f^{T|X=x}(\vartheta) = \frac{f^{X|T=\vartheta}(x)f^T(\vartheta)}{\int_{\Theta} f^{X|T=t}(x)f^T(t)\nu(dt)}, \quad \vartheta \in \Theta, \quad (1.4)$$

für  $\tilde{\mathbb{P}}^X$ -f.a.  $x \in \mathcal{X}$  mit dem Wahrscheinlichkeitsmaß  $\tilde{\mathbb{P}}$  aus Bemerkung 1.27. Das **a-posteriori-Risiko** eines Schätzers  $\hat{\rho}$ , gegeben  $X = x$ , ist definiert durch

$$R_{\pi}(\hat{\rho}|x) = \int_{\Theta} \ell(\vartheta, \hat{\rho}(x))f^{T|X=x}(\vartheta)\nu(d\vartheta). \quad (1.5)$$

Der Name der a-posteriori-Verteilung oder englisch *posterior distribution* suggeriert bereits, dass diese Verteilung nach einer gemachten Beobachtung  $x \in \mathcal{X}$  berechnet wird.

*Bemerkung 1.30 (a-posteriori-Verteilung und -Risiko)*

- Beachte, dass im Nenner von (1.4) die Randdichte

$$f^X(x) = \int_{\Theta} f^{X|T=t}(x) f^T(t) \nu(dt) \quad (1.6)$$

von  $X$  bezüglich  $\mu$  in  $(\mathcal{X} \times \Theta, \mathcal{F} \otimes \mathcal{F}_{\Theta}, \tilde{\mathbb{P}})$  steht, sodass (1.4) für  $\tilde{\mathbb{P}}^X$ -f.a.  $x \in \mathcal{X}$  wohldefiniert ist.

- Im diskreten Fall wird das Integral in (1.5) und im Nenner von (1.4) zu einer Summe, sodass wir genau die klassische Bayes-Formel erhalten.

*Beispiel 1.31 (a-posteriori-Verteilung und -Risiko)* Sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\theta})_{\theta \in \Theta})$  ein von  $\mu$  dominiertes statistisches Modell. Setze  $\Theta = \{0, 1\}$ ,  $\mathcal{F}_{\Theta} := \mathcal{P}(\Theta)$ ,  $\ell(\theta, r) = |\theta - r|$  (0-1-Verlust) und betrachte eine a-priori-Verteilung  $\pi$  mit  $\pi(\{0\}) =: \pi_0$  und  $\pi(\{1\}) =: \pi_1 = 1 - \pi_0$ . Die Wahrscheinlichkeitsmaße  $\mathbb{P}_0$  und  $\mathbb{P}_1$  mögen Dichten  $p_0$  und  $p_1$  bezüglich eines Maßes  $\mu$  besitzen (zum Beispiel  $\mu = \mathbb{P}_0 + \mathbb{P}_1$ ). Dann ist die a-posteriori-Verteilung auf  $\Theta$  durch die Zähldichte

$$f^{T|X=x}(i) = \frac{\pi_i p_i(x)}{\pi_0 p_0(x) + \pi_1 p_1(x)}, \quad i = 0, 1, \quad (\tilde{\mathbb{P}}^X\text{-f.ü.})$$

gegeben. Damit ist das a-posteriori-Risiko eines Schätzers  $\hat{\vartheta}: \mathcal{X} \rightarrow \{0, 1\}$  als Erwartungswert bezüglich der a-posteriori-Dichte gegeben durch

$$R_{\pi}(\hat{\vartheta}|x) = \frac{\hat{\vartheta}(x) \pi_0 p_0(x) + (1 - \hat{\vartheta}(x)) \pi_1 p_1(x)}{\pi_0 p_0(x) + \pi_1 p_1(x)}.$$

Das a-posteriori-Risiko ist minimal für  $\hat{\vartheta}_{\pi}(x) := \mathbb{1}(\pi_1 p_1(x) \geq \pi_0 p_0(x))$ , was genau dem später zu besprechenden *Bayes-Klassifizierer* entspricht (siehe Kapitel ??).

**Kurzbiografie (Thomas Bayes)** Thomas Bayes wurde um 1702 in London geboren. Er schrieb sich 1719 an der University of Edinburgh ein und studierte Logik und Theologie. Von etwa 1734 bis 1752 war er Pfarrer der presbyterianischen Kirche in Tunbridge Wells bei London. Erst in seinen späten Jahren vertiefte er sein Interesse an der Beschreibung von Wahrscheinlichkeiten und schrieb seine Ideen und Ergebnisse in Manuskripten nieder, die nach seinem Tod veröffentlicht wurden. Er war damit einer der ersten Statistiker, der sich mit Wahrscheinlichkeit befasste. Nicht nur der *Satz von Bayes* wurde nach ihm benannt, auch das große Gebiet der *Bayes-Statistik*. 1761 starb Thomas Bayes in Tunbridge Wells.

**Satz 1.32 (Bayes-Risiko und Bayes-Schätzer)** *Es gelten die Bedingungen der vorangegangenen Definition. Für das Bayes-Risiko eines Schätzers  $\hat{\rho}$  mit  $R_{\pi}(\hat{\rho}) < \infty$  gilt dann mit der Randdichte aus (1.6)*

$$R_\pi(\widehat{\rho}) = \int_{\mathcal{X}} R_\pi(\widehat{\rho}|x) f^X(x) \mu(dx).$$

Minimiert  $\widehat{\rho}(x)$  für  $\widetilde{\mathbb{P}}^X$ -fast alle  $x$  das a-posteriori-Risiko im Sinne von

$$R_\pi(\widehat{\rho}|x) = \min_{r \in \mathbb{R}^p} \int_{\Theta} \ell(\vartheta, r) f^{T|X=x}(\vartheta) \nu(d\vartheta),$$

dann ist  $\widehat{\rho}$  ein Bayes-Schätzer.

**Beweis** Umstellen von (1.4) ergibt  $f^{T|X=x}(\vartheta) f^X(x) = f^{X|T=\vartheta}(x) f^T(\vartheta)$ . Setzen wir dies für die Dichte von  $\mathbb{P}_\vartheta(dx) \pi(d\vartheta)$  ein, ergibt sich

$$\begin{aligned} R_\pi(\widehat{\rho}) &= \int_{\Theta} \int_{\mathcal{X}} \ell(\vartheta, \widehat{\rho}(x)) \mathbb{P}_\vartheta(dx) \pi(d\vartheta) \\ &= \int_{\Theta} \int_{\mathcal{X}} \ell(\vartheta, \widehat{\rho}(x)) f^{T|X=x}(\vartheta) f^X(x) \mu(dx) \nu(d\vartheta) \\ &= \int_{\mathcal{X}} R_\pi(\widehat{\rho}|x) f^X(x) \mu(dx), \end{aligned}$$

wobei wir im letzten Schritt den Satz von Fubini anwenden können, da die Integranden nichtnegativ sind. Minimiert  $\widehat{\rho}(x)$  punktweise den Integranden, so ist auch das Integral minimal.  $\square$

Dieser Satz liefert uns eine neue Methode zur Konstruktion von Schätzern.

**Korollar 1.33** *Unter quadratischem Verlust ist der Bayes-Schätzer gegeben durch den a-posteriori-Mittelwert*

$$\widehat{\rho}(x) = \int_{\Theta} \rho(\vartheta) f^{T|X=x}(\vartheta) \nu(d\vartheta) = \mathbb{E}[\rho(T)|X = x].$$

Wir verwenden hier den über bedingte Dichten definierten bedingten Erwartungswert  $\mathbb{E}[\rho(T)|X = x]$  für  $\widetilde{\mathbb{P}}^X$ -f.a.  $x \in \mathcal{X}$ , der sich natürlich in das abstraktere Konzept bedingter Erwartungen einbettet. Den Beweis dieses Korollars überlassen wir der Leserin als Übung. Für den wichtigen Fall, dass  $\rho(\vartheta) = \vartheta \in \mathbb{R}^p$  die Identität ist, ergibt sich der a-posteriori-Mittelwert

$$\widehat{\vartheta}(x) = \int_{\Theta} \vartheta f^{T|X=x}(\vartheta) \nu(d\vartheta) = \mathbb{E}[T|X = x]$$

als Bayes-Schätzer zum quadratischen Verlust. Oft ist  $\nu$  das Zählmaß oder das Lebesgue-Maß, sodass sich das vorangegangene Integral als Summe bzw. als Riemann-Integral schreiben lässt.

Analog ist der Bayes-Schätzer bezüglich absolutem Verlust für  $\rho(\vartheta) = \vartheta \in \mathbb{R}$  gegeben durch den *Median* der a-posteriori-Verteilung:

$$\widetilde{\vartheta}(x) = \inf \left\{ q \in \mathbb{R} : \int_{-\infty}^q f^{T|X=x}(\vartheta) \nu(d\vartheta) \geq \frac{1}{2} \right\}$$

Falls  $\nu$  ein Zählmaß ist und damit die a-posteriori-Verteilung diskret, ist auch der 0-1-Verlust von Interesse. Als Bayes-Schätzer ergibt sich dann die Maximalstelle  $\arg \max_{\vartheta \in \Theta} f^{T|X=x}(\vartheta)$  der a-posteriori-Verteilung, die man auch MAP-Schätzer (Maximum-a-posteriori-Schätzer) nennt. Motiviert durch die Maximum-Likelihood-Methode wird der Maximum-a-posteriori-Schätzer auch über den diskreten Fall hinaus eingesetzt. Es handelt sich dann allerdings nicht mehr um einen Bayes-Schätzer im Sinne von Satz 1.32.

**Methode 1.34 (Bayes-Schätzer)** Durch die Wahl einer a-priori-Verteilung ergeben sich durch den a-posteriori-Mittelwert, den a-posteriori-Median und die a-posteriori-Maximalstelle Schätzer für den unbekannten Parameter.

*Beispiel 1.35 (Bayes-Schätzer – a)* Wir betrachten das statistische Modell  $(\{0, \dots, n\}, \mathcal{P}(\{0, \dots, n\}), (\text{Bin}(n, \vartheta))_{\vartheta \in [0,1]})$  und wählen wie in Beispiel 1.28 als a-priori-Verteilung  $\pi = \text{U}([0, 1])$  mit Dichte  $f^T = \mathbb{1}_{[0,1]}$ . Um die a-posteriori-Verteilung zu berechnen, genügt es, deren Dichte bis auf die Normierungskonstante zu bestimmen, da wir bereits wissen, dass sich wieder eine Wahrscheinlichkeitsdichte in  $\vartheta$  ergibt. Wir verwenden das Symbol  $\propto$  für Gleichheit bis auf eine multiplikative Konstante. Es gilt für  $x \in \{0, \dots, n\}$  und  $\vartheta \in [0, 1]$

$$\begin{aligned} f^{T|X=x}(\vartheta) &= \frac{f^{X|T=\vartheta}(x)f^T(\vartheta)}{f^X(x)} \propto f^{X|T=\vartheta}(x)f^T(\vartheta) \\ &= \binom{n}{x} \vartheta^x (1-\vartheta)^{n-x} \\ &\propto \vartheta^x (1-\vartheta)^{n-x}. \end{aligned}$$

Wir erkennen die Struktur einer Beta-Verteilung  $\text{Beta}(\alpha, \beta)$  mit Parametern  $\alpha = x+1$  und  $\beta = n-x+1$ . Als Bayes-Schätzer unter quadratischem Verlust ergibt sich der Erwartungswert der a-posteriori-Verteilung

$$\hat{\vartheta} = \int_0^1 \vartheta f^{T|X=x}(\vartheta) d\vartheta = \frac{\alpha}{\alpha + \beta} = \frac{x+1}{n+2}.$$

Wir erhalten eine Alternative zum Maximum-Likelihood-Schätzer, die uns bereits in Beispiel 1.14 begegnet ist.

*Beispiel 1.36 (Bayes-Schätzer – b)* Sei  $X_1, \dots, X_n \sim \text{N}(\mu, \sigma^2)$  eine mathematische Stichprobe mit bekanntem  $\sigma^2 > 0$  und a-priori-Verteilung  $\pi = \text{N}(a, b^2)$ . Mithilfe der Bayes-Formel kann die a-posteriori-Verteilung für eine Realisierung  $x = (x_1, \dots, x_n)$  berechnet werden. Da  $f^{T|X=x}$  wieder eine Dichte, also normiert, sein muss, brauchen wir die Konstanten nicht mitzuberechnen und verwenden das Symbol  $\propto$  für Gleichheit bis auf eine multiplikative Konstante:

$$\begin{aligned} f^{T|X=x}(\mu) &\propto f^{X|T=\mu}(x)f^T(\mu) \\ &\propto \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right) \exp\left(-\frac{(\mu - a)^2}{2b^2}\right) \end{aligned}$$

$$\begin{aligned}
&\propto \exp\left(-\frac{\mu^2 - 2\mu\bar{x}_n}{2\sigma^2/n} - \frac{\mu^2 - 2a\mu}{2b^2}\right) \\
&= \exp\left(-\frac{(b^2 + \sigma^2/n)\mu^2 - 2\mu(b^2\bar{x}_n + a\sigma^2/n)}{2b^2\sigma^2/n}\right) \\
&\propto \exp\left(-\frac{1}{2}\left(\frac{n}{\sigma^2} + \frac{1}{b^2}\right)\left(\mu - \frac{b^2\bar{x}_n}{b^2 + \sigma^2/n} - \frac{a\sigma^2/n}{b^2 + \sigma^2/n}\right)^2\right)
\end{aligned}$$

Gegeben die Beobachtung  $X$ , ist  $\mu$  also gemäß

$$N\left(\frac{b^2}{b^2 + \sigma^2/n}\bar{X}_n + \frac{\sigma^2/n}{b^2 + \sigma^2/n}a, (n\sigma^{-2} + b^{-2})^{-1}\right)$$

a-posteriori-verteilt. Da hier der a-posteriori-Mittelwert mit dem a-posteriori-Median und der Maximalstelle der a-posteriori-Dichte übereinstimmt, ist der Bayes-Schätzer bezüglich quadratischem Verlust und absolutem Verlust sowie der MAP-Schätzer gegeben durch

$$\hat{\mu}_n = \frac{b^2}{b^2 + \sigma^2/n}\bar{X}_n + \frac{\sigma^2/n}{b^2 + \sigma^2/n}a.$$

Wir gewichten also den empirischen Mittelwert  $\bar{X}_n$  und den a-priori-Mittelwert  $a$  entsprechend dem Verhältnis der Varianzen  $\frac{\sigma^2}{n}$  und  $b^2$ .

In beiden Beispielen konnten wir die a-posteriori-Verteilung explizit bestimmen. Hierbei handelt es sich allerdings um Spezialfälle in denen die Berechnung von Bayes-Schätzern besonders einfach ist. Für komplexere Modelle führt die Berechnung der a-posteriori-Dichte mitunter auf hochdimensionale Integrale. Dazu werden häufig numerische Methoden wie MCMC-Methoden (Markov Chain Monte Carlo) verwendet.

**Bemerkung 1.37** Eine Familie von Verteilungen  $\mathcal{D}$  auf  $(\Theta, \mathcal{F}_\Theta)$  heißt durch  $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$  *konjugiert* (lateinisch für „verbunden“), falls für jede a-priori-Verteilung  $\pi \in \mathcal{D}$  und jede Beobachtung  $X = x$  die a-posteriori-Verteilung ebenfalls zu  $\mathcal{D}$  gehört. Beispiel 1.36 zeigt, dass die Verteilungsklasse der Normalverteilungen durch Normalverteilungen mit unbekanntem Erwartungswert und bekannter Varianz konjugiert werden. Beispiel 1.35 lässt sich auf Beta-verteilte a-priori-Verteilungen verallgemeinern, sodass auch hier eine konjugierte Verteilungsklasse vorliegt (Aufgabe 1.6).

Zum Schluss dieses Kapitels wollen wir noch einen Zusammenhang zwischen Minimax- und Bayes-Schätzer deutlich machen.

**Lemma 1.38** In einem dominierten statistischen Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit messbarem Parameterraum  $(\Theta, \mathcal{F}_\Theta)$  gilt für jeden Schätzer  $\hat{\rho}$

$$\sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho}) = \sup_{\pi} R_{\pi}(\hat{\rho}),$$

wobei sich das zweite Supremum über alle a-priori-Verteilungen  $\pi$  auf  $(\Theta, \mathcal{F}_\Theta)$  mit  $(\mathcal{F} \otimes \mathcal{F}_\Theta)$ -messbarer gemeinsamer Dichte  $f^{X|T=}$  erstreckt. Insbesondere ist das Risiko eines Bayes-Schätzer stets kleiner oder gleich dem Minimax-Risiko.

**Beweis** Natürlich gilt  $R_\pi(\hat{\rho}) = \int_{\Theta} R(\vartheta, \hat{\rho}) \pi(d\vartheta) \leq \sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho})$  für alle  $\pi$ . Andererseits folgt durch Betrachtung der a-priori-Punktverteilungen  $\delta_\vartheta$ ,  $\vartheta \in \Theta$ , dass  $\sup_\pi R_\pi(\hat{\rho}) \geq \sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho})$  gilt. Dies zeigt die Behauptung.  $\square$

Durch dieses Lemma erhalten wir untere Schranken für das Minimax-Risiko über das Risiko von Bayes-Schätzern. Mögliche Anwendungen illustriert der folgende Satz.

**Satz 1.39** Sei  $X_1, \dots, X_n$  eine  $N(\mu, \sigma^2)$ -verteilte mathematische Stichprobe mit unbekanntem Erwartungswert  $\mu \in \mathbb{R}$  und bekanntem  $\sigma^2 > 0$ . Bezüglich quadratischem Risiko ist das arithmetische Mittel  $\bar{X}_n$  ein Minimax-Schätzer von  $\mu$ .

**Beweis** Wir betrachten a-priori-Verteilungen  $\pi = N(0, b^2)$  für  $\mu$ . Nach Beispiel 1.36 ist die a-posteriori-Verteilung gegeben durch

$$N\left(\frac{b^2 \bar{X}_n}{b^2 + \sigma^2/n}, (n\sigma^{-2} + b^{-2})^{-1}\right).$$

Der Bayes-Schätzer bezüglich quadratischem Risiko ist gegeben durch den a-posteriori-Erwartungswert  $\hat{\mu}_n = b^2 \bar{X}_n / (b^2 + \sigma^2 n^{-1})$  und sein a-posteriori-Risiko ist gegeben durch die Varianz der a-posteriori-Verteilung. Ist  $f^X$  die Randdichte von  $X$  bezüglich  $\tilde{\mathbb{P}}$ , folgt aus Satz 1.32:

$$\begin{aligned} R_\pi(\hat{\mu}_n) &= \int_{\mathbb{R}^n} \text{Var}_{T|X=x}(\mu) f^X(x) dx \\ &= \int_{\mathbb{R}^n} (n\sigma^{-2} + b^{-2})^{-1} f^X(x) dx = (n\sigma^{-2} + b^{-2})^{-1} \end{aligned}$$

Somit können wir das Minimax-Risiko nach unten abschätzen:

$$\begin{aligned} \inf_{\tilde{\mu}} \sup_{\mu \in \mathbb{R}} R(\mu, \tilde{\mu}) &= \inf_{\tilde{\mu}} \sup_{\pi} R_\pi(\tilde{\mu}) \geq \inf_{\tilde{\mu}} \sup_{b>0} R_{N(0, b^2)}(\tilde{\mu}) \\ &\geq \sup_{b>0} \inf_{\tilde{\mu}} R_{N(0, b^2)}(\tilde{\mu}) \\ &= \sup_{b>0} (n\sigma^{-2} + b^{-2})^{-1} = \frac{\sigma^2}{n}, \end{aligned}$$

wie behauptet, da  $R(\mu, \bar{X}_n) = \sigma^2/n$ .  $\square$

Damit ist  $\hat{\mu}_1$  aus Beispiel 1.24 minimax.

**Bemerkung 1.40** Auch in höheren Dimensionen ist das koeffizientenweise Stichprobenmittel  $\bar{X}_n$  ein Minimax-Schätzer für den Mittelwertvektor der Normalverteilung, das heißt das maximale Risiko ist kleinstmöglich. Ab Dimension 3 gibt es jedoch Schätzer, insbesondere den sogenannten *James-Stein-Schätzer*, deren Risiko für jedes  $\mu \in \mathbb{R}$  stets kleiner ist als das Risiko von  $\bar{X}_n$ . Dieser Effekt ist als *Stein-Phänomen* bekannt, siehe Lehmann and Casella (1998). In der Sprache der Entscheidungstheorie wird  $\bar{X}_n$  als *nicht zulässig* bezeichnet. Hierbei gibt es keinen Widerspruch zur Minimax-Eigenschaft von  $\bar{X}_n$ , denn die Verbesserung des James-Stein-Schätzers wird beliebig klein für  $|\mu| \rightarrow \infty$ .

### 1.3 Hypothesentests

Wir wenden uns nun der zweiten grundlegenden Fragestellung zu. Häufig ist weniger die gesamte zugrunde liegende Verteilung von Interesse als die Frage, ob eine bestimmte Eigenschaft erfüllt ist. Hierfür formuliert man die gefragte Eigenschaft als Hypothese und entscheidet dann, ob die gemachten Beobachtungen für oder gegen diese Hypothese sprechen. Die Entscheidung dafür oder dagegen beruht auf einem statistischen (Hypothesen-)Test.

#### 1.3.1 Statistische Tests und ihre Fehler

Als einführendes Beispiel kann man sich ein Gerichtsverfahren vorstellen, bei dem eine Person beschuldigt wird, eine Straftat begangen zu haben. Das Gericht geht prinzipiell von der Annahme aus, dass der Angeklagte unschuldig ist. Diese Annahme nennen die Juristen Unschuldsvermutung, wir Statistiker nennen sie abstrakter *Nullhypothese*  $H_0$ . Die *Alternativhypothese*  $H_1$  lautet, dass der Angeklagte schuldig ist. Es wird nun versucht, die Schuld des Angeklagten zu erörtern.

Das Gerichtsverfahren soll verhindern, den Angeklagten zu Unrecht zu verurteilen. Die Wahrscheinlichkeit für diesen sogenannten *Fehler 1. Art* soll anhand von Indizien und Argumentationen möglichst gering sein. Dies beeinflusst aber auch den sogenannten *Fehler 2. Art*: Ein Schuldiger sollte nicht freigesprochen werden. Beide Fehlerarten werden nicht gleichzeitig minimiert werden können. Die Unschuldsvermutung impliziert eine asymmetrische Gewichtung der Fehler erster und zweiter Art. Ganz ähnlich werden wir bei der Konstruktion statistischer Tests vorgehen.

**Definition 1.41** Wir betrachten ein statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ . Die Parametermenge sei in zwei disjunkte Teilmengen  $\Theta_0$  und  $\Theta_1$  zerlegt, das heißt  $\Theta = \Theta_0 \cup \Theta_1$  und  $\Theta_0 \cap \Theta_1 = \emptyset$ . Das **Testproblem** liest sich dann als

$$H_0: \vartheta \in \Theta_0 \quad \text{gegen} \quad H_1: \vartheta \in \Theta_1.$$

Dabei werden  $H_0, H_1$  als **Hypothesen** bezeichnet, genauer heißt  $H_0$  **Nullhypothese** und  $H_1$  **Alternativhypothese** oder **Alternative**.

*Beispiel 1.42 (Hypothesentest)* Denken wir wieder an unser Saskia-Beispiel 1.7. Die Studierenden in der Umfrage hatten nur zwei Auswahlmöglichkeiten: Sie stimmten für 1, wenn sie mit ihrem Studium zufrieden waren, und für 0, wenn Sie mit ihrem Studium unzufrieden waren. Die Antworten  $X_1, \dots, X_n$  hatten wir als Bernoulli-verteilt modelliert mit Parameter  $\vartheta = \mathbb{P}(X_i = 1) \in (0, 1)$ .

1. Laut einer Umfrage unter 6000 Studierenden<sup>1</sup> sind 74,9% der Studierendenschaft eher zufrieden bis ganz und gar zufrieden mit ihrem Studium. Dieses Ergebnis kann Saskia anhand ihrer eigenen Erhebung überprüfen. Die entsprechende Nullhypothese lautet dann formal  $H_0 : \vartheta \in [\frac{3}{4}, 1) =: \Theta_0$ . Da wir die Parametermenge  $\Theta = (0, 1)$  nach obiger Definition in disjunkte Teilmengen aufteilen müssen, bleibt für die Alternativhypothese  $H_1$  der Parameterbereich  $\Theta_1 = (0, \frac{3}{4})$  übrig. Das Testproblem lautet ausformuliert

$$H_0 : \vartheta \in [\frac{3}{4}, 1) \quad \text{gegen} \quad H_1 : \vartheta \in (0, \frac{3}{4}).$$

2. Saskia könnte auch die Hypothese aufstellen, dass das Verhältnis zwischen Studierenden, die mit ihrem Studium zufrieden bzw. unzufrieden sind, ausgewogen ist. Dies würde zur Nullhypothese  $H_0 : \vartheta = 1/2$  und Alternative  $H_1 : \vartheta \neq 1/2$  führen.

Ein statistischer Test entscheidet nun zwischen Nullhypothese und Alternative aufgrund einer Beobachtung  $x \in \mathcal{X}$ .

**Definition 1.43** Ein **nichtrandomisierter statistischer Test** ist eine messbare Abbildung

$$\varphi : (\mathcal{X}, \mathcal{F}) \rightarrow (\{0, 1\}, \mathcal{P}(\{0, 1\}))$$

mit dem Testentscheid

$$\varphi(x) = \begin{cases} 1 : & \text{die Nullhypothese wird verworfen/abgelehnt,} \\ 0 : & \text{die Nullhypothese wird nicht verworfen bzw. wird akzeptiert.} \end{cases}$$

Die Menge  $\{\varphi = 1\} = \{x \in \mathcal{X} : \varphi(x) = 1\}$  heißt **Ablehnbereich** oder **kritischer Bereich** von  $\varphi$ .

Ein **randomisierter statistischer Test** ist eine messbare Abbildung  $\varphi : (\mathcal{X}, \mathcal{F}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$ . Im Fall  $\varphi(x) \in (0, 1)$  entscheidet ein unabhängiges Bernoulli-Zufallsexperiment mit Erfolgswahrscheinlichkeit  $p = \varphi(x)$ , ob die Nullhypothese verworfen wird.

Testen beinhaltet mögliche Fehlentscheidungen:

- (i) **Fehler 1. Art** (auch  $\alpha$ -Fehler, englisch: *type I error* oder *false positive* genannt): Entscheidung für  $H_1$ , obwohl  $H_0$  wahr ist,
- (ii) **Fehler 2. Art** (auch  $\beta$ -Fehler, englisch: *type II error* oder *false negative* genannt): Entscheidung für  $H_0$ , obwohl  $H_1$  wahr ist.

Ein nicht-randomisierter Test lässt sich als randomisierter Test mit Werten nur in  $\{0, 1\}$  interpretieren. Wenn wir allgemein von Tests sprechen, meinen wir daher formal randomisierte Tests. In den Beispielen werden wir uns zunächst auf nicht-randomisierte (statistische) Tests konzentrieren. Im folgenden Abschnitt wird sich

<sup>1</sup> HIS (2008). *Studenten, die gegenwärtig mit den folgenden Bereichen ihres Lebens eher zufrieden bis ganz und gar zufrieden sind*. In Statista: <https://de.statista.com/statistik/daten/studie/1440/umfrage/zufriedenheit-von-studenten/> (Zugegriffen am 05. November 2020).

jedoch auch zeigen, dass randomisierte Tests ihre (mathematische) Berechtigung haben.

Im einführenden Gerichtsbeispiel hatten wir die Schwierigkeit angemerkt, eine sinnvolle Schwelle zu finden, ab wann genügend Indizien und Argumente für eine Verurteilung vorliegen. Wir wollen nicht, dass ein Unschuldiger zu Unrecht verurteilt wird (Fehler 1. Art), und wir wollen auch nicht, dass ein Schuldiger frei gesprochen wird (Fehler 2. Art). In dieser Analogie soll also ein statistischer Test entscheiden, ob genügend Indizien und Argumente vorliegen.

**Definition 1.44** Sei  $\varphi$  ein Test der Hypothese  $H_0: \vartheta \in \Theta_0$  gegen die Alternative  $H_1: \vartheta \in \Theta_1$  im statistischen Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ . Die **Gütefunktion** von  $\varphi$  ist definiert als

$$\beta_\varphi: \Theta \rightarrow \mathbb{R}_+, \vartheta \mapsto \mathbb{E}_\vartheta[\varphi].$$

Ein Test  $\varphi$  erfüllt das **Signifikanzniveau**  $\alpha \in [0, 1]$  (oder  $\varphi$  ist ein **Test zum Niveau**  $\alpha$ ), falls  $\beta_\varphi(\vartheta) \leq \alpha$  für alle  $\vartheta \in \Theta_0$  gilt.

Erfüllt ein Test ein *vorher* festgelegtes Niveau  $\alpha$ , so nennen wir das Ergebnis des Tests statistisch *signifikant* (zum Niveau  $\alpha$ ). In Definition 1.44 legen wir also  $\alpha \in [0, 1]$  vorher fest und konstruieren den Test  $\varphi$  anschließend so, dass für alle Parameter  $\vartheta \in \Theta_0$  der Nullhypothese der Erwartungswert  $\mathbb{E}_\vartheta[\varphi]$  kleiner gleich  $\alpha$  ist.

Ist  $\varphi: (\mathcal{X}, \mathcal{F}) \rightarrow (\{0, 1\}, \mathcal{P}(\{0, 1\}))$  ein nicht-randomisierter Test, so vereinfacht sich diese Forderung: Da

$$\mathbb{E}_\vartheta[\varphi] = 0 \cdot \mathbb{P}_\vartheta(\varphi = 0) + 1 \cdot \mathbb{P}_\vartheta(\varphi = 1) = \mathbb{P}_\vartheta(\varphi = 1)$$

gilt, erfüllt  $\varphi$  das Signifikanzniveau  $\alpha$ , falls

$$\beta_\varphi(\vartheta) = \mathbb{P}_\vartheta(\varphi = 1) \leq \alpha \quad \text{für alle } \vartheta \in \Theta_0$$

gilt. Die Wahrscheinlichkeit, dass unser Test die Nullhypothese ablehnt, obwohl  $H_0$  stimmt, soll also höchstens  $\alpha$  sein. Das Signifikanzniveau  $\alpha$  beschränkt damit die Wahrscheinlichkeit für Fehler 1. Art.

Greifen wir noch einmal die Analogie zum Gerichtswesen auf. Für Juristen ist es schlimmer, Unschuldige zu Unrecht zu verurteilen als Schuldige fälschlicherweise freizusprechen (Unschuldsvermutung). Für sie sind also Fehler 1. Art schwerwiegender als Fehler 2. Art. Deshalb wird versucht, die Fehler 1. Art durch ein geringes Signifikanzniveau  $\alpha$  klein zu halten. Typische Werte sind  $\alpha = 0,05$  und  $\alpha = 0,01$ .

Auf der Alternativmenge  $\Theta_1$  sollte die Gütefunktion möglichst groß sein, damit der Fehler 2. Art klein ist. Für einen nichtrandomisierten Test  $\varphi$  ergibt sich direkt der Zusammenhang

$$\mathbb{P}_\vartheta(\varphi = 0) = 1 - \mathbb{P}_\vartheta(\varphi = 1) = 1 - \beta_\varphi(\vartheta) \quad \text{für } \vartheta \in \Theta_1.$$

Entsprechend gibt  $1 - \beta_\varphi(\vartheta)$  auch für randomisierte Tests die Wahrscheinlichkeit eines Fehler 2. Art an.

In der Regel lassen sich die Wahrscheinlichkeiten für Fehler 1. und 2. Art nicht gleichzeitig minimieren. Deshalb verfährt man im Allgemeinen so, dass

- (i) zuerst die Wahrscheinlichkeit für Fehler 1. Art durch ein vorgegebenes Signifikanzniveau begrenzt wird und dann
- (ii) unter der Maßgabe von (i) ein Test  $\varphi$  gesucht wird, der die Wahrscheinlichkeit für Fehler 2. Art minimiert.

Als Nullhypothese  $H_0$  sollte im Allgemeinen der bisherige Standard oder die aufgrund theoretischer Überlegungen erzielte Modellierung gewählt werden. Wird  $H_0$  nicht zugunsten von  $H_1$  abgelehnt, heißt das noch lange nicht, dass  $H_0$  „wahr“ ist, sondern nur, dass die Daten nicht im Widerspruch zur Nullhypothese stehen und  $H_1$  die Daten nicht besser als  $H_0$  erklärt. Wird hingegen  $H_0$  abgelehnt, besteht Grund zu der Annahme, dass die Modelle unter  $H_0$  die Daten nicht hinreichend gut beschreiben.

Als Nächstes betrachten wir zwei Formen von Testproblemen, die oft in der Praxis auftauchen.

**Definition 1.45** Es sei  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell, wobei  $\Theta \subseteq \mathbb{R}$  ein zusammenhängendes Intervall ist.

1. Testprobleme der Form  $H_0: \vartheta \leq \vartheta_0$  gegen  $H_1: \vartheta > \vartheta_0$  oder  $H_0: \vartheta \geq \vartheta_0$  gegen  $H_1: \vartheta < \vartheta_0$  für ein  $\vartheta_0 \in \Theta$  heißen **einseitig**.
2. Testprobleme der Form  $H_0: \vartheta = \vartheta_0$  gegen  $H_1: \vartheta \neq \vartheta_0$  für ein  $\vartheta_0 \in \Theta$  heißen **zweiseitig**.

Hypothesentests für einseitige bzw. zweiseitige Testprobleme werden selbst **einseitig** bzw. **zweiseitig** genannt.

*Beispiel 1.46 (Einseitiger Binomialtest)* Von den 13 Todesfällen unter 55- bis 65-jährigen Arbeitern eines Kernkraftwerks im Jahr 1995 waren 5 auf einen Tumor zurückzuführen. Die Todesursachenstatistik 1995 weist aus, dass Tumore bei etwa 1/5 aller Todesfälle die Ursache in der betreffenden Altersklasse (in der Gesamtbevölkerung) sind. Ist die beobachtete Häufung von tumorbedingten Todesfällen signifikant zum Niveau 5%? Das heißt, ist die Wahrscheinlichkeit dafür, dass 5 von 13 Kernkraftwerksarbeitern an Tumoren gestorben sind, in Anbetracht dessen, dass allgemein jeder 5. an Tumoren stirbt, kleiner oder gleich 5%?

Um diese Frage zu beantworten, beschreiben wir die Anzahl der Tumortoten als Zufallsvariable  $X \in \{0, 1, \dots, n\}$  mit  $n = 13$ . Als statistisches Modell, in das  $X$  abbildet, wählen wir  $\mathcal{X} = \{0, \dots, n\}$ ,  $\mathcal{F} = \mathcal{P}(\mathcal{X})$  und  $\mathbb{P}_p = \text{Bin}(13, p)$ . Der Parameter  $p \in [0, 1]$  ist die Wahrscheinlichkeit, dass eine Person an einem Tumor gestorben ist. Wir wollen wissen, ob 5 von 13 Todesfällen eine signifikante Häufung zum Niveau  $\alpha = 5\%$  sind. Diese Fragestellung führt auf das Testproblem

$$H_0: p \leq 1/5 \quad \text{gegen} \quad H_1: p > 1/5.$$

Wir müssen nun einen geeigneten (nichtrandomisierten) Test  $\varphi$  zum Niveau  $\alpha = 0,05$  konstruieren. Naheliegenderweise wählen wir

$$\varphi(x) := \mathbb{1}(x > c). \tag{1.7}$$

Hierbei wollen wir den kritischen Wert  $c > 0$  so wählen, dass  $\beta_\varphi(p) \leq \alpha$  für alle  $p \leq 1/5$  gilt. Wegen der zugrunde liegenden Binomialverteilung nennen wir diesen Test *Binomialtest*. Für  $p \leq 1/5$  gilt

$$\beta_\varphi(p) = \mathbb{P}_p(\varphi = 1) = \mathbb{P}_p(X > c) \leq \sup_{p \leq 1/5} \mathbb{P}_p(X > c) \stackrel{!}{\leq} \alpha. \quad (1.8)$$

Um eine möglichst große Güte des Tests zu erreichen, sollte  $c$  unter dieser Nebenbedingung möglichst klein gewählt werden. Für die Verteilungsfunktion  $\mathbb{P}_p(X \leq k)$  unserer Binomialverteilung gilt für  $k \in \mathcal{X}$

$$\mathbb{P}_p(X \leq k) = \sum_{l=0}^k \binom{13}{l} p^l (1-p)^{13-l}.$$

Da  $p \mapsto \mathbb{P}_p(X \leq k)$  für alle  $k \in \mathcal{X}$  monoton fallend auf  $[0, 1]$  ist (dies sieht man durch Ableiten), folgt, dass  $p \mapsto \mathbb{P}_p(X > c) = 1 - \mathbb{P}_p(X \leq c)$  monoton steigend ist. Folglich gilt mit (1.8), dass

$$\sup_{p \leq 1/5} \mathbb{P}_p(X > c) = \mathbb{P}_{1/5}(X > c) \stackrel{!}{\leq} \alpha = 0,05 \quad \Leftrightarrow \quad \mathbb{P}_{1/5}(X \leq c) \stackrel{!}{\geq} 0,95.$$

Da  $c$  eine natürliche Zahl sein sollte und wegen

$$\mathbb{P}_{1/5}(X \leq 4) \approx 0,901 \quad \text{und} \quad \mathbb{P}_{1/5}(X \leq 5) \approx 0,970,$$

wählen wir  $c = 5$ , siehe Abbildung 1.4 auf Seite 34. Unser Test lautet also  $\varphi(x) = \mathbb{1}_{\{x > 5\}}$ . Setzen wir nun die Anzahl 5 von Kernkraftwerksarbeitern, die an einem Tumor gestorben sind, in  $\varphi$  ein, so erhalten wir 0. Das heißt, unser Test  $\varphi$  akzeptiert zum Niveau  $\alpha = 0,05$  die Nullhypothese, dass die beobachteten Todesfälle nicht signifikant sind. Wir erhalten also

$$\sup_{p \leq 1/5} \beta_\varphi(p) = \beta_\varphi(1/5) \approx 0,03 < 0,05 = \alpha,$$

und der Binomialtest kann aufgrund des diskreten Stichprobenraums das Niveau nicht voll ausschöpfen.

Dieses Beispiel führt uns auf ein allgemeines Konstruktionsprinzip von Tests einer Hypothese  $H_0: \vartheta \in \Theta_0$  gegen die Alternative  $H_1: \vartheta \in \Theta_1$  mit  $\Theta_0 \neq \emptyset$  und  $\Theta_1 = \Theta \setminus \Theta_0$ .

**Methode 1.47 (Konstruktion von Tests durch kritische Werte)** Für ein gegebenes Testproblem und ein festgelegtes Signifikanzniveau  $\alpha \in (0, 1)$  wähle eine **Teststatistik** bzw. **Test-** oder **Prüfgröße**  $T: (\mathcal{X}, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  und konstruiere einen Niveau- $\alpha$ -Test durch

$$\varphi(x) = \mathbb{1}(T(x) > c), \quad x \in \mathcal{X}, \quad (1.9)$$

wobei  $c \in \mathbb{R}$  so gewählt wird, dass

$$\sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T(X) > c) \leq \alpha \quad (1.10)$$

gilt. Das kleinstmögliche  $c = c_{\alpha}$ , sodass (1.10) erfüllt ist, heißt **kritischer Wert**. Analog sind auch kritische Bereiche der Formen  $\{T(x) \geq c\}$ ,  $\{T(x) < c\}$ ,  $\{T(x) \leq c\}$ ,  $\{|T(x)| > c\}$ ,  $\{T(x) < u\} \cup \{T(x) > o\}$  usw. möglich.

Für eine möglichst hohe Güte des Tests möchten wir den Ablehnbereich möglichst groß wählen und minimieren daher  $c$  unter der Bedingung (1.10). Um den kritischen Wert zu wählen und somit das Signifikanzniveau des Tests zu gewährleisten, muss die Verteilung der Teststatistik unter der Nullhypothese bekannt sein. Da man häufig Monotonieeigenschaften nutzen kann (wie im Beispiel 1.46), vereinfacht sich die Wahl von  $c_{\alpha}$  im Fall von ein- oder zweiseitigen Testproblemen oft, wenn das Gleichheitszeichen in der Nullhypothese steht (siehe Definition 1.45). Im Spezialfall  $\Theta_0 = \{\vartheta_0\}$  ist für einen einseitigen Test der kritische Wert  $c_{\alpha}$  genau das  $(1 - \alpha)$ -Quantil der Verteilung von  $T$  unter  $\mathbb{P}_{\vartheta_0}$ .

Natürlich hängt sowohl die Teststatistik als auch der kritische Bereich vom jeweiligen Testproblem ab. Die Herausforderung ist also die Wahl einer Teststatistik, die (gut) geeignet ist, um die jeweilige Hypothese zu überprüfen. Nachdem wir in Beispiel 1.46 bereits einen einseitigen Binomialtest gesehen haben, illustriert das folgende Beispiel den zweiseitigen Fall.

*Beispiel 1.48 (Zweiseitiger Binomialtest)* Wir wollen die Hypothese: „Es werden genauso viele Jungen wie Mädchen geboren.“ überprüfen.

Sind bei  $n \in \mathbb{N}$  Geburten  $w$  Mädchen zur Welt gekommen, ist es sinnvoll, als Stichprobenraum  $\mathcal{X} = \{0, \dots, n\}$  zu wählen und als statistisches Modell  $(\mathcal{X}, \mathcal{P}(\mathcal{X}), (\mathbb{P}_{\vartheta})_{\vartheta \in [0,1]})$  mit Binomialverteilungen  $\mathbb{P}_{\vartheta} = \text{Bin}(n, \vartheta)$ . Die Hypothese führt auf das zweiseitige Testproblem

$$H_0 : \vartheta = \vartheta_0 \quad \text{gegen} \quad H_1 : \vartheta \neq \vartheta_0$$

für  $\vartheta_0 = \frac{1}{2}$ , wobei  $w \in \mathcal{X}$  beobachtet wird. Wir setzen das Niveau  $\alpha = 0,05$ . Die Teststatistik  $T(w) = w$  führt auf den *zweiseitigen Binomialtest*

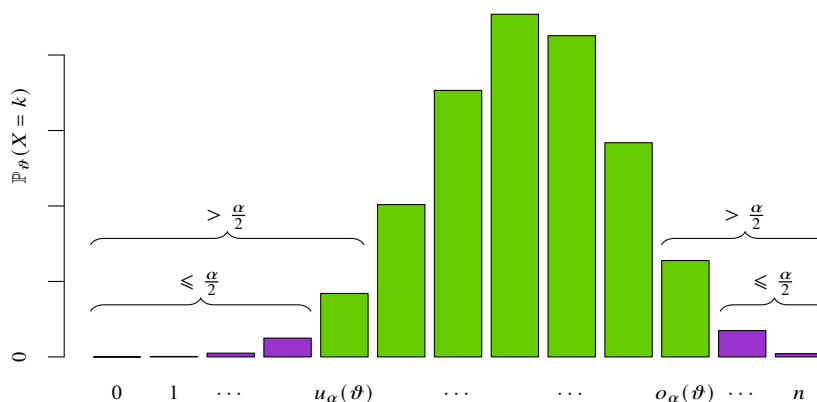
$$\varphi_{\alpha}(w) = 1 - \mathbb{1}_{[u_{\alpha}(\vartheta_0), o_{\alpha}(\vartheta_0)]}(w)$$

mit kritischen Werten

$$\begin{aligned} u_{\alpha}(\vartheta_0) &:= \max\{k \in \mathbb{N} : \mathbb{P}_{\vartheta_0}(\{0, \dots, k-1\}) \leq \alpha/2\} \quad \text{und} \\ o_{\alpha}(\vartheta_0) &:= \min\{k \in \mathbb{N} : \mathbb{P}_{\vartheta_0}(\{k+1, \dots, n\}) \leq \alpha/2\}. \end{aligned} \quad (1.11)$$

Die symmetrische Verteilung der Fehlerwahrscheinlichkeiten am unteren und am oberen Rand ist dabei eine vernünftige Wahl, wobei auch jede andere Aufteilung des Niveaus  $\alpha$  prinzipiell möglich ist.

Beachte, dass dieser zweiseitige Binomialtest nur im Fall  $\vartheta_0 = 1/2$  symmetrisch um  $\vartheta_0$  ist, denn nur dann ist die Binomialverteilung symmetrisch, siehe



**Abb. 1.2** Obere und untere Grenze des Annahmebereichs des zweiseitigen Binomialtests für ein  $\vartheta > 1/2$

Abbildung 1.2. Im symmetrischen Fall erhalten wir die einfachere Darstellung  $\varphi_\alpha(w) = 1 - \mathbb{1}(|\frac{w}{n} - \frac{1}{2}| \leq c_\alpha)$  mit einem geeigneten kritischen Wert  $c_\alpha$ .

Laut Statistischem Bundesamt wurden im Jahr 2018 in Hamburg 21.126 Kinder (lebend) geboren. Durch die Berechnung der Quantile der Binomialverteilung führt dies auf

$$u_{0,05}(1/2) = 10.421 \quad \text{und} \quad o_{0,05}(1/2) = 10.705 \quad \text{bzw.} \quad c_{0,05} = 0,00672.$$

Von diesen 21.126 Kindern waren 10.215 Mädchen. Setzen wir dies in  $\varphi_\alpha$  ein, erhalten wir

$$\varphi_{0,05}(10.215) = 1 - \mathbb{1}\left(\left|\frac{10.215}{21.126} - \frac{1}{2}\right| \leq c_{0,05}\right) = 1 - 0 = 1,$$

sodass unser Test  $\varphi_\alpha$  zum Niveau  $\alpha = 0,05$  die Nullhypothese „Es werden genauso viele Jungen wie Mädchen geboren.“ ablehnt.

*Bemerkung 1.49 (Normalapproximation)* Aufgrund des zentralen Grenzwertsatzes kann die Binomialverteilung  $\text{Bin}(n, \vartheta)$  bei hinreichend großen Stichprobenumfängen durch eine Normalverteilung approximiert werden. Es bietet sich in diesem Fall also an, einen *Gauß-Test* zu verwenden, um den Binomialtest zu approximieren: Für  $\vartheta \in (0, 1)$  normalisieren wir die Beobachtung  $X \sim \text{Bin}(n, \vartheta)$  durch

$$Y := \frac{X - n\vartheta}{\sqrt{n\vartheta(1-\vartheta)}}.$$

Aus dem zentralen Grenzwertsatz folgt dann, dass die Verteilung von  $Y$  für  $n \rightarrow \infty$  gegen  $N(0, 1)$  konvergiert. Für die Teststatistik  $T(X) := |X/n - \vartheta|$  und eine standardnormalverteilte Zufallsvariable  $Z \sim N(0, 1)$  erhalten wir

$$\begin{aligned}
\mathbb{P}_\vartheta(T(X) > c_\alpha) &= \mathbb{P}_\vartheta\left(\frac{|X - n\vartheta|}{\sqrt{n\vartheta(1-\vartheta)}} > \sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right) \\
&\stackrel{n \rightarrow \infty}{\approx} \mathbb{P}\left(|Z| > \sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right) \\
&= 2\left(1 - \mathbb{P}\left(Z \leq \sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right)\right) \\
&= 2\left(1 - \Phi\left(\sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right)\right) \stackrel{!}{=} \alpha,
\end{aligned}$$

wobei  $\Phi$  die Verteilungsfunktion der Standardnormalverteilung bezeichnet. Ist die Gleichheit für einen kritischen Wert  $c_\alpha$  erfüllt, erhalten wir einen nichtrandomisierten Test, der asymptotisch (!) für  $n \rightarrow \infty$  das Niveau  $\alpha$  erreicht. Umformen ergibt

$$c_\alpha = \sqrt{\frac{\vartheta_0(1-\vartheta_0)}{n}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

mit  $\vartheta = \vartheta_0$  unter  $H_0$ .

Die Normalapproximation hat den Vorteil, dass im Vergleich zur diskreten Binomialverteilung die Berechnung deutlich vereinfacht wird. Aufgrund der heutigen leistungsfähigen Computer ist dieses Argument in der Praxis allerdings zu vernachlässigen. Zusätzlich vereinfacht die Approximation jedoch die Interpretation der Teststatistik und ermöglicht einen Vergleich zwischen Studien mit verschiedenen großen Stichprobenumfängen.

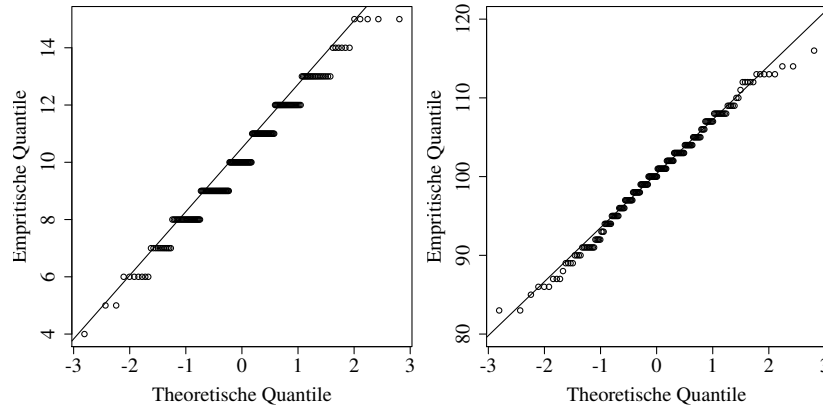
Ob die Normalapproximation der Binomialverteilung tatsächlich passend ist, wird in der Praxis mit sogenannten QQ-Plots überprüft.

*Bemerkung 1.50 (QQ-Plots)* Ein *Quantil-Quantil-Plot*, *QQ-Plot* oder auch *QQ-Diagramm* ist ein exploratives, grafisches Werkzeug, in dem die Quantile zweier Zufallsvariablen gegeneinander aufgetragen werden, um ihre Verteilungen zu vergleichen. Da sie im Allgemeinen von großer Bedeutung ist, erklären wir im Detail, wie man die Verteilung einer Beobachtung mit der Standardnormalverteilung vergleicht: Die *empirische Verteilungsfunktion* einer mathematischen Stichprobe  $X_1, \dots, X_N$  ist definiert als

$$F_N(x) := \frac{1}{N} \sum_{i=1}^N \mathbb{1}(X_i \leq x).$$

Für große  $N$  approximiert  $F_N$  die wahre Verteilungsfunktion  $F$ , da nach dem starken Gesetz der großen Zahlen  $F_N(x) \rightarrow \mathbb{E}[\mathbb{1}(X_1 \leq x)] = F(x)$   $\mathbb{P}$ -f.s. für alle  $x \in \mathbb{R}$  gilt (tatsächlich gilt diese Konvergenz sogar gleichmäßig auf  $\mathbb{R}$  nach dem Satz von Glivenko-Cantelli). Im Fall  $X_i \sim N(\mu, \sigma^2)$  gilt  $F(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$ . Für die Quantilfunktion ergibt sich also

$$F^{-1}(\Phi(x)) = \Phi^{-1}(\Phi(x)) \cdot \sigma + \mu = \sigma \cdot x + \mu,$$



**Abb. 1.3** QQ-Plot für eine binomialverteilte Stichprobe mit den Parameter  $p = 1/2$ ,  $n = 20$  (links) sowie  $n = 200$  (rechts)

und  $F^{-1} \circ \Phi$  beschreibt eine Gerade. Im QQ-Plot werden der Größe nach geordnete Werte  $(x_k)$  aus dem Intervall  $[0, 1]$  in die verallgemeinerte Inverse  $F_n^{-1}$  und in  $\Phi^{-1}$  eingesetzt, sodass man für jedes  $x_k$  ein Quantilpaar erhält. Diese Paare werden (als Koordinatenpaare interpretiert) in ein Koordinatensystem eingetragen, und wenn die  $X_i \sim N(\mu, \sigma^2)$  sind, sollten die Quantilpaare in etwa auf einer Geraden liegen.

Zur Illustration betrachten wir eine Stichprobe  $X_1, \dots, X_{200} \sim \text{Bin}(n, 1/2)$  für  $n \in \{20, 200\}$  in Abbildung 1.3. Die QQ-Plots zeigen deutlich, dass die Normalapproximation der Binomialverteilung für kleine  $n$  problematisch ist, aber für große  $n$  gut funktioniert. Genauer sollte  $np(1-p)$  hinreichend groß sein, was man aus dem Satz von Berry-Esseen folgern kann. Abbildung 1.3 macht allerdings auch deutlich, dass selbst für  $n = 200$  die extremen Quantile der Binomialverteilung noch deutlich von der Normalapproximation abweichen.

*Beispiel 1.51 (Planung des Stichprobenumfangs)* In Beispiel 1.1 hatten wir uns gefragt, wie viele Studierende an der Umfrage von Saskia teilnehmen sollten, um ein aussagekräftiges Ergebnis zu erhalten. Wie wir im Folgenden sehen werden, kann diese Frage durch eine Betrachtung des Fehlers 2. Art beantwortet werden.

Gemäß Beispiel 1.7 wählen wir das statistische Modell

$$(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\text{Ber}(\vartheta))^{\otimes n}_{\vartheta \in [0, 1]}).$$

Die Summe  $Y = X_1 + \dots + X_n$  der zufriedenen Studentinnen und Studenten ist also wieder binomialverteilt. In Beispiel 1.42 interessierten wir uns für das Testproblem  $H_0 : \vartheta \geq \vartheta_0$  gegen  $H_1 : \vartheta < \vartheta_0$  mit  $\vartheta_0 = 3/4$ . Wir legen ein Signifikanzniveau  $\alpha \in (0, 1)$  fest und betrachten den Test  $\varphi_n(x) = \mathbb{1}_{\{x_1 + \dots + x_n < c_n\}}$ . Eine Normalapproximation liefert den kritischen Wert

$$c_n = n\vartheta_0 + \Phi^{-1}(\alpha)\sqrt{n\vartheta_0(1-\vartheta_0)}$$

(dem/der Lesenden sei die Überprüfung überlassen, dass dies tatsächlich ein Test zum asymptotischen Niveau  $\alpha$  ist). Wir modifizieren nun die Alternative etwas zu  $H_1 : \vartheta < \vartheta_1$  für ein  $\vartheta_1 < \vartheta_0$ , sodass die Parameter der Hypothese und der Alternative voneinander getrennt sind. Wir wollen die Wahrscheinlichkeit für Fehler 2. Art auf höchstens  $1 - \beta$  begrenzen, das heißt, es soll  $\inf_{\vartheta \leq \vartheta_1} \beta_{\varphi_n}(\vartheta) = \beta_{\varphi_n}(\vartheta_1) \geq \beta$  gelten (die Gütefunktion ist monoton fallend). Wir berechnen

$$\begin{aligned}
 \beta &\stackrel{!}{\leq} \beta_{\varphi_n}(\vartheta_1) \\
 &= \mathbb{P}_{\vartheta_1} \left( \sum_{j=1}^n X_j < c_n \right) \\
 &= \mathbb{P}_{\vartheta_1} \left( \frac{\sum_{j=1}^n X_j - n\vartheta_1}{\sqrt{n\vartheta_1(1-\vartheta_1)}} < \frac{c_n - n\vartheta_1}{\sqrt{n\vartheta_1(1-\vartheta_1)}} \right) \\
 &= \mathbb{P}_{\vartheta_1} \left( \frac{\sum_{j=1}^n X_j - n\vartheta_1}{\sqrt{n\vartheta_1(1-\vartheta_1)}} < \frac{\sqrt{n}(\vartheta_0 - \vartheta_1) + \Phi^{-1}(\alpha)\sqrt{\vartheta_0(1-\vartheta_0)}}{\sqrt{\vartheta_1(1-\vartheta_1)}} \right) \\
 &\stackrel{ZGWS}{\approx} \Phi \left( \sqrt{n} \frac{\vartheta_0 - \vartheta_1}{\sqrt{\vartheta_1(1-\vartheta_1)}} + \Phi^{-1}(\alpha) \sqrt{\frac{\vartheta_0(1-\vartheta_0)}{\vartheta_1(1-\vartheta_1)}} \right).
 \end{aligned}$$

Wir erhalten damit die Bedingung

$$n \geq \frac{\vartheta_1(1-\vartheta_1)}{(\vartheta_0 - \vartheta_1)^2} \left( \Phi^{-1}(\beta) - \Phi^{-1}(\alpha) \sqrt{\frac{\vartheta_0(1-\vartheta_0)}{\vartheta_1(1-\vartheta_1)}} \right)^2.$$

Insbesondere sehen wir, dass mehr Beobachtungen benötigt werden, um den Fehler 2. Art auch nahe der Nullhypothese, das heißt wenn  $\vartheta_0 - \vartheta_1$  klein ist, durch  $1 - \beta$  zu beschränken.

Setzen wir  $\vartheta_0 = 3/4$ ,  $\vartheta_1 = 0,7$ ,  $\alpha = 0,1$  und  $\beta = 0,9$ , dann erhalten wir  $n \approx 522$ , was relativ wenig ist bei einer Grundgesamtheit von 20.000 Studierenden. Wollen wir hingegen die Wahrscheinlichkeiten für Fehler 1. und 2. Art weiter reduzieren, sagen wir mit  $\alpha = 0,05$  und  $\beta = 0,95$ , erhalten wir  $n \approx 860$ .  $\alpha = 0,01$  und  $\beta = 0,99$  führen zu  $n \approx 1720$ .

Ein wichtiges Konzept in der Anwendung statistischer Tests sind *p*-Werte. Um diese zu motivieren, bleiben wir beim Saskia-Beispiel. Wir wollen aber das Testproblem etwas umformulieren. Die Nullhypothese ist, dass die Studierenden eher sehr zufrieden mit ihrem Studium sind ( $H_0 : \vartheta \in [3/4, 1)$ ), und wir wollen testen, ob sie mit ihrem Studium eher nicht zufrieden sind ( $H_1 : \vartheta \in (0, 3/4)$ ). Nehmen wir an, wir erhalten durch die Umfrage eine Stichprobe, die im arithmetischen Mittel 0,5 ergibt. Die Hälfte ist also mit dem Studium (un-)zufrieden. Nun könnte man sich fragen, wie stark dieses Ergebnis der Nullhypothese widerspricht, denn immerhin ist die andere Hälfte der Stichprobe zufrieden und es haben (mit sehr hoher Wahr-

scheinlichkeit) nicht alle Studierenden die Umfrage beantwortet. Eine Antwort auf diese Frage liefert der *p-Wert*.

**Definition 1.52** Sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $\varphi = \mathbb{1}(T > c_\alpha)$  ein Test der Hypothese  $H_0: \vartheta \in \Theta_0 \neq \emptyset$  mit Teststatistik  $T$  und kritischem Wert  $c_\alpha$ . Dann ist der **p-Wert**  $p_\varphi(x)$  einer Realisierung  $x \in \mathcal{X}$  definiert als

$$p_\varphi(x) = \inf_{t < t^*} \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) > t) \quad \text{für } t^* = T(x).$$

Der p-Wert ist also die Wahrscheinlichkeit, bei Gültigkeit der Hypothese etwas mindestens so Extremes zu beobachten, wie das tatsächlich Beobachtete. Um den Zusammenhang zu Niveau- $\alpha$ -Tests zu erkennen, betrachten wir zunächst den Spezialfall einer einfachen Hypothese  $\Theta_0 = \{\vartheta_0\}$  und einer stetigen, streng monoton wachsenden Verteilungsfunktion von  $T$  unter  $\mathbb{P}_{\vartheta_0}$ . Dann vereinfacht sich obige Definition zu

$$p_\varphi(x) = \mathbb{P}_{\vartheta_0}(T(X) > T(x)).$$

Unter diesen Bedingungen gilt zudem  $\mathbb{P}_{\vartheta_0}(T > c_\alpha) = \alpha$ . Wir können nun durch einen Vergleich des p-Werts mit dem Niveau ablesen, ob die Hypothese abgelehnt wird oder nicht:

- Falls  $p_\varphi(x) < \alpha$ , dann folgt aus der Monotonie von  $c \mapsto \mathbb{P}_{\vartheta_0}(T > c)$ , dass  $T(x) > c_\alpha$ . Der Test  $\varphi(x) = \mathbb{1}(T(x) > c_\alpha)$  lehnt die also Hypothese ab.
- Falls  $p_\varphi(x) > \alpha$ , ergibt sich analog  $T(x) < c_\alpha$ , sodass der Test die Hypothese nicht verwirft.

Wir werden diesen Zusammenhang auch unter deutlich allgemeineren Bedingungen in Satz 1.53 nachweisen. Statt nur zu prüfen, ob ein Test eine Hypothese akzeptiert oder ablehnt, gibt der p-Wert (auch „Signifikanzwahrscheinlichkeit“, „Überschreitungswahrscheinlichkeit“ oder „Signifikanzwert“) das kleinste Signifikanzniveau an, zu dem eine Hypothese abgelehnt würde. Damit gibt der p-Wert Aufschluss darüber, „wie stark“ die Daten der Hypothese widersprechen.

Der p-Wert spielt in der Wissenschaft (Biologie, Chemie, Medizin, Physik, Soziologie, ...) eine wichtige Rolle als Qualitätsmaß für die Relevanz beziehungsweise Richtigkeit einer Theorie. Zum Beispiel wurde in der Arbeit Abbott et al. (2016), die erstmals die Existenz der Gravitationswellen nachgewiesen hat, ein p-Wert von  $7,5 \times 10^{-8}$  angegeben. Zudem hat der p-Wert für die Theorie des *multiplen Testens* eine große Bedeutung und ist als Zwischenresultat für Methoden, die beispielsweise die *false discovery rate* kontrollieren, eine fundamentale Größe, siehe Dickhaus (2014).

**Satz 1.53 (Eigenschaften des p-Werts)** Es seien  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell,  $\alpha_0 \in (0, 1)$  und  $\varphi = \mathbb{1}(T > c_{\alpha_0})$  ein Niveau- $\alpha_0$ -Test der Hypothese  $H_0: \vartheta \in \Theta_0 \neq \emptyset$  mit einer Teststatistik  $T: \mathcal{X} \rightarrow \mathbb{R}$  und kritischen Werten

$$c_\alpha = \inf \left\{ c \in \mathbb{R} : \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) > c) \leq \alpha \right\}, \quad \alpha \in (0, 1). \quad (1.12)$$

Unter der Annahme

$$\sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T > c_{\alpha}) \leq \alpha \quad \text{für alle } \alpha \in (0, 1)$$

gilt die Darstellung

$$p_{\varphi}(x) := \inf_{\alpha: T(x) > c_{\alpha}} \sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T(X) > c_{\alpha}). \quad (1.13)$$

sowie der Testentscheid

$$\varphi(x) = \begin{cases} 1, & \text{falls } p_{\varphi}(x) < \alpha_0, \\ 0, & \text{falls } p_{\varphi}(x) > \alpha_0. \end{cases}$$

**Bemerkung 1.54 (Kritische Werte)** Im Allgemeinen muss die Eigenschaft  $\sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T(X) > c_{\alpha}) \leq \alpha$  für die kritischen Werte  $c_{\alpha}$  aus (1.12) nicht erfüllt sein. Dieses Defizit kann behoben werden indem man die Gültigkeit der Vertauschungsbedingung

$$\lim_{c_n \downarrow c} \sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T > c_n) = \sup_{\vartheta \in \Theta_0} \lim_{c_n \downarrow c} \mathbb{P}_{\vartheta}(T > c_n) \quad (1.14)$$

fordert. In diesem Fall folgt aus der  $\sigma$ -Stetigkeit von Maßen

$$\sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T > c_{\alpha}) = \sup_{\vartheta \in \Theta_0} \lim_{c_n \downarrow c_{\alpha}} \mathbb{P}_{\vartheta}(T > c_n) = \lim_{c_n \downarrow c_{\alpha}} \sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}(T > c_n) \leq \alpha.$$

Gleichung (1.14) gilt beispielsweise für einfache Hypothesen, das heißt für  $\Theta_0 = \{\vartheta_0\}$ , da dann  $\sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta} = \mathbb{P}_{\vartheta_0}$  ein Maß ist, oder wenn Monotonie der Verteilungsfunktion von  $T(X)$  in  $\vartheta$  wie in Beispiel 1.56 unten verwendet werden kann.

**Beweis** Im Folgenden schreiben wir  $\mathbb{P}_0 := \sup_{\vartheta \in \Theta_0} \mathbb{P}_{\vartheta}$ . Wir beginnen mit folgenden elementaren Beobachtungen:

- (i) Aus der Monotonie der Maße  $\mathbb{P}_{\vartheta}$  folgt sofort, dass die Abbildung  $c \mapsto \mathbb{P}_0(T > c)$  monoton fallend ist.
- (ii) Sei  $c < t^*$ ,  $\alpha = \mathbb{P}_0(T > c)$  und  $c_{\alpha}$  der entsprechende kritische Wert. Aus der Definition der kritischen Werte folgt sofort  $c_{\alpha} \leq c$ , und laut Voraussetzung gilt  $\mathbb{P}_0(T > c_{\alpha}) \leq \alpha$ . Aus der Monotonie (i) folgt dann aber

$$\alpha = \mathbb{P}_0(T > c) \leq \mathbb{P}_0(T > c_{\alpha}) \leq \alpha,$$

und somit  $\mathbb{P}_0(T > c_{\alpha}) = \alpha$ .

Setze  $S := \{c_{\alpha} < t^*, \alpha \in [0, 1]\}$  die Menge aller kritischen Werte kleiner als  $t^*$ . Dann gilt  $S \neq \emptyset$ , denn für ein beliebiges  $c < t^*$  und den kritischen Wert  $c_{\alpha}$  zu  $\alpha = \mathbb{P}_0(T > c)$  gilt  $c_{\alpha} \leq c < t^*$  und somit  $c_{\alpha} \in S$ . Wir betrachten nun das Supremum  $s := \sup\{c_{\alpha} < t^*\} \leq t^*$  der kritischen Werte kleiner  $t^*$  und unterscheiden die Fälle  $s = t^*$  und  $s < t^*$ .

Falls  $s = t^*$  gilt, gibt es eine Folge kritischer Werte  $(c_{\alpha,n})$  mit  $\lim_{n \rightarrow \infty} c_{\alpha,n} = t^*$  und wir erhalten direkt die Darstellung

$$p_\varphi(x) = \inf_{t < t^*} \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) > t) = \inf_{\alpha: t^* > c_\alpha} \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) > c_\alpha).$$

Wir betrachten nun den Fall  $s < t^*$ . Wir zeigen zuerst, dass  $s$  ein kritischer Wert zum Niveau  $\alpha_s = \mathbb{P}_0(T > s)$  ist. Wegen der Monotonie (i) und (ii) gilt für alle  $c_\alpha \in S$

$$\mathbb{P}_0(T > c_\alpha) \geq \mathbb{P}_0(T > s) = \alpha_s = \mathbb{P}_0(T > c_{\alpha_s}).$$

Aus der Monotonie (i) folgt dann  $c_{\alpha_s} \geq c_\alpha$  für alle  $c_\alpha \in S$ . Damit ist  $c_{\alpha_s}$  eine obere Schranke für  $S$  und muss aufgrund von  $c_{\alpha_s} \leq s$  mit dem Supremum  $s$  übereinstimmen. Aus  $c_{\alpha_s} = s < t^*$  und der Monotonie (i) folgt daher

$$\mathbb{P}_0(T > c_{\alpha_s}) \geq \inf_{\alpha: t^* > c_\alpha} \mathbb{P}_0(T > c_\alpha) \geq \mathbb{P}_0(T > s), \quad (1.15)$$

d.h.  $\inf_{\alpha: t^* > c_\alpha} \mathbb{P}_0(T > c_\alpha) = \mathbb{P}_0(T > s)$ .

Als Nächstes zeigen wir, dass zwischen  $s$  und  $t^*$  keine „Masse“ mehr ist. Dazu argumentieren wir per Widerspruch: Angenommen, es gibt ein  $s < s' < t^*$ , sodass  $\mathbb{P}_0(s' < T) < \mathbb{P}_0(s < T) = \alpha_s$ . Für  $\alpha' = \mathbb{P}_0(s' < T)$  mit kritischem Wert  $c_{\alpha'} \leq s'$  folgt dann

$$\alpha_s = \mathbb{P}_0(T > s) > \mathbb{P}_0(T > s') = \alpha' \stackrel{(ii)}{=} \mathbb{P}_0(T > c_{\alpha'}).$$

Insbesondere muss als  $c_{\alpha'} > s$  gelten im Widerspruch zu  $c_{\alpha'} \leq s'$ . Damit ergibt sich

$$\inf_{s \leq s' < t^*} \mathbb{P}_0(T > s') = \mathbb{P}_0(T > s),$$

und zusammen mit (1.15) folgt die Gültigkeit von (1.13).

Man beachte nun, dass  $\varphi(x) = 1$  genau dann gilt, wenn  $T(x) = t^* > c_{\alpha_0}$ . Die Annahme an die kritischen Werte und die Monotonie (i) von  $\mathbb{P}_0$  liefern nun

$$\alpha_0 \geq \mathbb{P}_0(T > c_{\alpha_0}) \geq \inf_{\alpha: c_\alpha < t^*} \mathbb{P}_0(T > c_\alpha) = p_\varphi(x).$$

Äquivalent impliziert  $p_\varphi(x) > \alpha_0$ , dass  $\varphi(x) = 0$  gelten muss. Im Fall  $p_\varphi(x) < \alpha_0$  folgt aus

$$p_\varphi(x) = \inf_{c < t^*} \mathbb{P}_0(T > c),$$

dass es ein  $c < t^*$  mit  $p_\varphi(x) \leq \mathbb{P}_0(T > c) \leq \alpha_0$  gibt. Die Definition der kritischen Werte impliziert nun, dass  $c \geq c_{\alpha_0}$  gelten muss, und wegen  $t^* > c \geq c_{\alpha_0}$  folgt  $\varphi(x) = 1$ .  $\square$

#### Bemerkung 1.55 ( $p$ -Wert)

1. Unter den Voraussetzungen von Satz 1.53 kann man den  $p$ -Wert alternativ auch durch

$$p_\varphi(x) = \inf\{\alpha \in [0, 1] : t^* > c_\alpha\}$$

darstellen (Übung 1.9).

2. Alle Rahmenbedingungen des Experiments, insbesondere also das Signifikanzniveau, müssen vor seiner Durchführung festgelegt werden! Ein Signifikanzniveau darf nicht a posteriori aufgrund der erzielten p-Werte bestimmt werden. Dies widerspricht korrekter statistischer Praxis! Mathematisch wäre  $\alpha$  eine Zufallsvariable (als Funktion in den Beobachtungen), und der vorangegangene Satz kann nicht angewendet werden.
3. Der Vorteil von p-Werten ist, dass sie unabhängig von einem a priori festgesetzten Signifikanzniveau  $\alpha$  berechnet werden können. Deshalb werden in allen gängigen Statistik-Softwaresystemen statistische Hypothesentests über die Berechnung von p-Werten implementiert.

*Beispiel 1.56 (p-Wert)* Auf Grundlage von i.i.d. normalverteilten Beobachtungen  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  mit unbekanntem Mittelwert  $\mu \in \mathbb{R}$  und bekanntem  $\sigma > 0$  soll das Testproblem

$$H_0: \mu \leq \mu_0 \quad \text{gegen} \quad H_1: \mu > \mu_0$$

untersucht werden. Da wir bereits wissen, dass im resultierenden statistischen Modell  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R})^{\otimes n}, (\mathbb{P}_\mu)_{\mu \in \mathbb{R}})$  mit  $\mathbb{P}_\mu = N(\mu, \sigma^2)^{\otimes n}$  das Stichprobenmittel  $\bar{X}_n$  ein guter Schätzer von  $\mu$  (sogar ein Minimax-Schätzer) ist, wählen wir den Test  $\varphi(x) = \mathbb{1}_{\{\bar{x}_n > c_\alpha\}}$  für kritische Werte  $(c_\alpha)_{\alpha \in (0,1)}$  und der Teststatistik  $T(x) = \bar{x}_n$ . Wegen  $\bar{X}_n \sim N(\mu, \sigma^2/n)$  unter  $\mathbb{P}_\mu$  und der Monotonie der Verteilungsfunktion  $\Phi$  der Standardnormalverteilung gilt für jedes  $c \in \mathbb{R}$

$$\sup_{\mu \leq \mu_0} \mathbb{P}_\mu(\bar{X}_n > c) = \sup_{\mu \leq \mu_0} \left\{ 1 - \Phi\left(\frac{\sqrt{n}(c - \mu)}{\sigma}\right) \right\} = 1 - \Phi\left(\frac{\sqrt{n}(c - \mu_0)}{\sigma}\right).$$

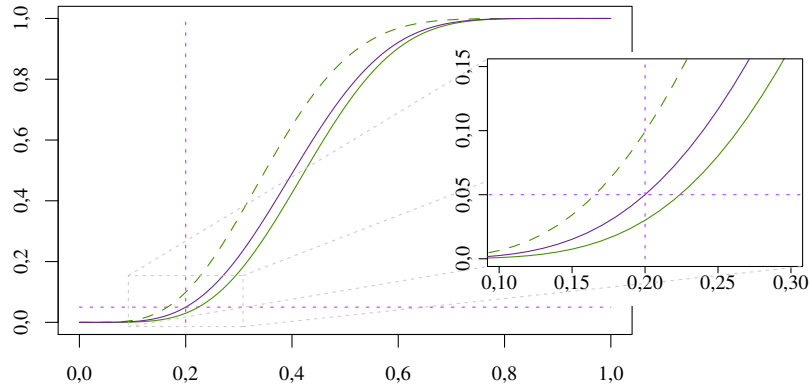
Zum einen folgt hieraus  $c_\alpha = \mu_0 + \frac{\sigma}{\sqrt{n}} q_{1-\alpha}$  für das  $(1 - \alpha)$ -Quantil von  $N(0, 1)$  und zum anderen ist die Bedingung  $\sup_{\mu \leq \mu_0} \mathbb{P}_\mu(\bar{X}_n > c_\alpha) \leq \alpha$  aus Satz 1.53 erfüllt. Wir erhalten für den *einseitigen Gauß-Test*  $\varphi(x) = \mathbb{1}_{\{x > \mu_0 + \frac{\sigma}{\sqrt{n}} q_{1-\alpha}\}}$  und für Realisierungen  $x \in \mathbb{R}^n$  die p-Werte

$$p_\varphi(x) = 1 - \Phi\left(\frac{\sqrt{n}(\bar{x}_n - \mu_0)}{\sigma}\right) = \inf\{\alpha \in (0, 1) : T(x) > c_\alpha\}.$$

### 1.3.2 Das Neyman-Pearson-Lemma

Wir werden nun ein Optimalitätskriterium für Tests kennenlernen. Zur Motivation kommen wir auf Beispiel 1.46 zurück.

*Beispiel 1.57 (Randomisierte Tests)* Wir betrachten erneut ein Binomialmodell  $\text{Bin}(13, p)$  auf  $\mathcal{X} = \{0, \dots, n\}$ . Da die Schwelle  $c$  des Binomialtests eine ganze



**Abb. 1.4** Die Gütefunktionen der Tests  $\varphi(x) := \mathbb{1}_{\{x>5\}}$  (grün, durchgezogen),  $\bar{\varphi}(x) := \mathbb{1}_{\{x>4\}}$  (grün, gestrichelt) sowie des randomisierten Tests  $\tilde{\varphi}$  (violett) im  $\text{Bin}(13, 1/5)$ -Modell. Violett markiert sind das Niveau  $\alpha = 0,05$  und der Parameterwert  $p = 1/5$ .

Zahl ist, kann der Test  $\varphi(x) := \mathbb{1}_{\{x>c\}}$  das Niveau  $\alpha = 0,05$  insofern nicht voll ausschöpfen, als dass die Wahrscheinlichkeit für Fehler 2. Art nicht perfekt minimiert werden kann. An dieser Stelle hilft eine Randomisierung des Tests. Wir betrachten einen Test der Form

$$\tilde{\varphi}(x) = \begin{cases} 0, & x < c, \\ \gamma, & x = c, \\ 1, & x > c. \end{cases}$$

Statt die Hypothese wie in Beispiel 1.46 bei  $X = c$  immer abzulehnen, fordern wir nun, dass im Fall  $X = c$  die Nullhypothese mit einer Wahrscheinlichkeit  $\gamma$  abgelehnt wird, und führen hierzu ein unabhängiges Bernoulli-Experiment  $B \sim \text{Ber}(\gamma)$  mit Erfolgswahrscheinlichkeit  $\gamma$  durch. Wir wählen  $\gamma$  so, dass das gesamte Signifikanzniveau ausgeschöpft wird, das heißt

$$0,05 \stackrel{!}{=} \mathbb{E}_{1/5}[\tilde{\varphi}] = \mathbb{P}_{1/5}(X > 5) + \gamma \cdot \mathbb{P}_{1/5}(X = 5).$$

Durch Umstellen nach  $\gamma$  ergibt sich  $\gamma \approx 0,29$ . Der resultierende randomisierte Test  $\tilde{\varphi}$  besitzt somit das Niveau  $\alpha = 0,05$  und ist zudem unverfälscht, siehe Abbildung 1.4. Da die Gütefunktion des Binomialtests stets unter der Gütefunktion des randomisierten Tests liegt, hat Letzterer einen kleineren Fehler zweiter Art.

In diesem Beispiel ist es uns gelungen, den Binomialtest zu verbessern. Es stellt sich die Frage, ob es auch einen besten Test gibt und, falls das der Fall ist, wie man diesen konstruieren kann. Die Antwort hierauf wollen wir im einfachen binären Modell studieren.

In diesem Fall besteht die Parametermenge nur aus zwei Parametern, sagen wir  $\Theta = \{0, 1\}$ . Wir möchten  $H_0: \vartheta = 0$  gegen  $H_1: \vartheta = 1$  testen. Das folgende grundlegende Resultat beschreibt Tests, die zu vorgegebenem Niveau die Wahrscheinlichkeit von Fehlern zweiter Art minimieren.

**Satz 1.58 (Neyman-Pearson-Lemma)** *Betrachte ein binäres statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \{0,1\}})$  mit Likelihood-Funktion  $L$  bezüglich eines dominierenden Maßes  $\mu$  (zum Beispiel  $\mu = \mathbb{P}_0 + \mathbb{P}_1$ ) sowie das Testproblem  $H_0: \vartheta = 0$  gegen  $H_1: \vartheta = 1$ . Dann besitzt der (möglicherweise randomisierte) Test*

$$\varphi_\alpha(x) := \begin{cases} 1, & \text{falls } L(1, x) > c_\alpha L(0, x), \\ 0, & \text{falls } L(1, x) < c_\alpha L(0, x), \\ \gamma_\alpha, & \text{falls } L(1, x) = c_\alpha L(0, x) \end{cases} \quad (1.16)$$

mit Konstanten  $c_\alpha \geq 0$ ,  $\gamma_\alpha \in [0, 1]$  unter allen (auch randomisierten) Tests  $\varphi$  mit demselben Niveau  $\mathbb{E}_0[\varphi] \leq \mathbb{E}_0[\varphi_\alpha]$  die kleinste Wahrscheinlichkeit für Fehler zweiter Art:

$$\mathbb{E}_1[1 - \varphi_\alpha] = \min_{\varphi} \mathbb{E}_1[1 - \varphi]$$

**Beweis** Das Resultat folgt, wenn wir  $E_1[\varphi - \varphi_\alpha] \leq 0$  gezeigt haben. Hierzu verwenden wir ein geschicktes Maßwechselargument sowie  $\varphi - \varphi_\alpha = \varphi - 1 \leq 0$  auf  $\{L(1) > c_\alpha L(0)\}$  und  $\varphi - \varphi_\alpha = \varphi \geq 0$  auf  $\{L(1) < c_\alpha L(0)\}$ :

$$\begin{aligned} \mathbb{E}_1[\varphi - \varphi_\alpha] &= \int_{\{L(1) > c_\alpha L(0)\}} (\varphi - \varphi_\alpha) L(1) d\mu + \int_{\{L(1) < c_\alpha L(0)\}} (\varphi - \varphi_\alpha) L(1) d\mu \\ &\quad + \int_{\{L(1) = c_\alpha L(0)\}} (\varphi - \varphi_\alpha) L(1) d\mu \\ &\leq \int_{\{L(1) > c_\alpha L(0)\}} (\varphi - \varphi_\alpha) c_\alpha L(0) d\mu \\ &\quad + \int_{\{L(1) < c_\alpha L(0)\}} (\varphi - \varphi_\alpha) c_\alpha L(0) d\mu \\ &\quad + \int_{\{L(1) = c_\alpha L(0)\}} (\varphi - \varphi_\alpha) c_\alpha L(0) d\mu \\ &= c_\alpha \mathbb{E}_0[\varphi - \varphi_\alpha] \\ &\leq 0, \end{aligned}$$

wobei wir in der letzten Zeile die Niveaubedingung  $\mathbb{E}_0[\varphi] \leq \mathbb{E}_0[\varphi_\alpha]$  verwendet haben.  $\square$

**Kurzbiografie (Jerzy Neyman)** Jerzy Neyman (bzw. Yuri Czeslawovich) wurde 1894 in Bendery im Russischen Kaiserreich geboren. Er studierte Mathematik und Physik in Charkiw (Ukraine). Nach dem Studium arbeitete er zunächst an den Universitäten in Charkiw und Warschau, bevor er 1934 von Egon Pearson ans University College London geholt wurde. 1938 bekam er einen Ruf an die University of

California in Berkeley, an der er das Statistik-Department aufbaute. 1981 starb Neymann in Kalifornien. Zu den wichtigsten Beiträgen Neymans zählen die Einführung von Konfidenzintervallen und das zusammen mit Egon Pearson bewiesene *Neyman-Pearson-Lemma*.

**Definition 1.59** In einem binären statistischen Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \{0,1\}})$  mit Likelihood-Funktion  $L$  bezüglich eines dominierenden Maßes heißt ein Test der Form (1.16) **Neyman-Pearson-Test**.

*Bemerkung 1.60*

1. Es sei der Leserin überlassen, zu zeigen, dass zu jedem Niveau  $\alpha \in (0, 1)$  durch Wahl von  $c_\alpha \geq 0$  und  $\gamma_\alpha \in [0, 1]$  ein Neyman-Pearson-Test  $\varphi_\alpha$  vom Niveau  $\alpha$  existiert (Aufgabe 2.10). Im Fall einer einelementigen Hypothese und Alternative, man spricht auch von *einfacher Hypothese* und *einfacher Alternative*, kann somit stets ein optimaler Test angegeben werden. Aus dem vorherigen Beweis lässt sich auch folgern, dass jeder in diesem Sinne optimale Test fast sicher von der Form eines Neyman-Pearson-Tests ist.
2. Die Neyman-Pearson-Tests hängen nicht von der Wahl des dominierenden Maßes  $\mu$  ab. Mittels Radon-Nikodym-Ableitungen sieht man nämlich mit dem kanonischen dominierenden Maß  $\bar{\mu} = \mathbb{P}_0 + \mathbb{P}_1$  und  $\bar{L}(\vartheta, x) = \frac{d\mathbb{P}_\vartheta}{d\bar{\mu}}(x)$ :

$$L(\vartheta, x) = \frac{d\mathbb{P}_\vartheta}{d\mu}(x) = \frac{d\mathbb{P}_\vartheta}{d\bar{\mu}}(x) \frac{d\bar{\mu}}{d\mu}(x) = \bar{L}(\vartheta, x) \frac{d\bar{\mu}}{d\mu}(x).$$

Die Likelihood-Funktionen  $L$  und  $\bar{L}$  unterscheiden sich also nur um einen in  $\vartheta$  konstanten Faktor, der sich bei der Konstruktion des Neyman-Pearson-Tests herauskürzt. Beachte dazu  $\mathbb{P}_\vartheta(\frac{d\bar{\mu}}{d\mu} = 0) = 0$  für  $\vartheta \in \{0, 1\}$  (Warum gilt das?).

*Beispiel 1.61 (Einseitiger Binomialtest)* Wir nehmen Beispiel 1.46 wieder auf, testen aber nun für  $n = 13$  und  $p_0 = 1/5$ ,  $p_1 = 1/4$  die einfache Hypothese „Tod durch Tumorerkrankung mit Wahrscheinlichkeit  $p = p_0$ “ gegen die einfache Alternative „Tod durch Tumorerkrankung mit Wahrscheinlichkeit  $p = p_1$ “. Wir erhalten das binäre statistische Modell  $(\{0, \dots, n\}, \mathcal{P}(\{0, \dots, n\}), (\mathbb{P}_\vartheta)_{\vartheta \in \{0,1\}})$  mit  $\mathbb{P}_0 = \text{Bin}(n; p_0)$  und  $\mathbb{P}_1 = \text{Bin}(n; p_1)$ . Mit dem Zählmaß als dominierendem Maß  $\mu$  erhalten wir den *Likelihood-Quotienten*

$$\frac{L(1, k)}{L(0, k)} = \frac{\binom{n}{k} p_1^k (1 - p_1)^{n-k}}{\binom{n}{k} p_0^k (1 - p_0)^{n-k}} = \frac{(1 - p_1)^n}{(1 - p_0)^n} \left( \frac{p_1(1 - p_0)}{p_0(1 - p_1)} \right)^k.$$

Wegen  $p_1 > p_0$  wächst dieser Quotient in  $k \in \{0, \dots, n\}$  streng monoton, und es gibt zu jedem  $c_\alpha \geq 0$  ein  $\tilde{c}_\alpha \geq 0$ , sodass

$$\varphi_\alpha(k) := \begin{cases} 1, & \text{falls } L(1, k) > c_\alpha L(0, k) \\ 0, & \text{falls } L(1, k) < c_\alpha L(0, k) \\ \gamma_\alpha, & \text{falls } L(1, k) = c_\alpha L(0, k) \end{cases} = \begin{cases} 1, & \text{falls } k > \tilde{c}_\alpha, \\ 0, & \text{falls } k < \tilde{c}_\alpha, \\ \gamma_\alpha, & \text{falls } k = \tilde{c}_\alpha. \end{cases}$$

Wählen wir nun wie in Beispiel 1.57  $\tilde{c}_{0,05} = 5$ ,  $\gamma_{0,05} \approx 0,29$ , so erreicht  $\varphi_{0,05}$  genau das Niveau  $\alpha = 0,05$  und besitzt als Neyman-Pearson-Test unter allen randomisierten Tests zum Niveau  $\alpha = 0,05$  die kleinste Fehlerwahrscheinlichkeit zweiter Art.

Der aufmerksamen Leserin ist vielleicht aufgefallen, dass die gesamte Testkonstruktion hier nicht vom exakten Wert  $p_1$  abhängt. Wir haben nur  $p_1 > p_0$  benutzt. Somit gilt sogar die weit stärkere Eigenschaft, dass es keinen Test  $\varphi$  von  $H_0: p = p_0$  gegen  $H_1: p > p_0$  vom Niveau  $\alpha = 0,05$  gibt, der für irgendein  $p \in (p_0, 1]$  eine größere Güte  $\mathbb{E}_p[\varphi]$  besitzt als  $\varphi_{0,05}$ . Wenn wir noch das Monotonieargument aus Beispiel 1.46 heranziehen, dass  $\varphi_{0,05}$  auch für die zusammengesetzte Hypothese  $H_0: p \leq p_0$  Niveau  $\alpha$  besitzt, so können wir weiter schließen, dass der einseitige Binomialtest auf der Alternative maximale Güte unter allen Niveau- $\alpha$ -Tests für  $H_0: p \leq p_0$  gegen  $H_1: p > p_0$  besitzt. Solche Tests nennt man auch *UMP-Tests* (*uniformly most powerful tests*).

**Beispiel 1.62 (Einseitiger Gauß-Test)** Wir beobachten eine mathematische Stichprobe normalverteilter Zufallsvariablen  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  mit  $\sigma > 0$  bekannt und wollen für feste  $\mu_0, \mu_1 \in \mathbb{R}$  die Hypothese  $H_0: \mu = \mu_0$  gegen  $H_1: \mu = \mu_1$  testen. Mit  $\mathbb{P}_0 = N(\mu_0, \sigma^2)^{\otimes n}$ ,  $\mathbb{P}_1 = N(\mu_1, \sigma^2)^{\otimes n}$  und dem  $n$ -dimensionalen Lebesgue-Maß als dominierendem Maß erhalten wir den Likelihood-Quotienten (als Statistik in den Beobachtungen geschrieben):

$$\begin{aligned} \frac{L(1)}{L(0)} &= \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n \left( (X_i - \mu_1)^2 - (X_i - \mu_0)^2 \right) \right) \\ &= \exp \left( -\frac{n}{\sigma^2} \left( (\mu_0 - \mu_1) \bar{X}_n + \frac{\mu_1^2}{2} - \frac{\mu_0^2}{2} \right) \right) \end{aligned}$$

mit dem Stichprobenmittel  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Betrachten wir von nun an den Fall  $\mu_1 > \mu_0$ , sehen wir, dass der Likelihood-Quotient eine streng monotone Funktion in  $\bar{X}_n$  ist und sonst nicht von den Beobachtungen abhängt. Also können wir wie beim Binomialtest durch Modifikation des kritischen Werts  $c_\alpha$  jeden Neyman-Pearson-Test schreiben als

$$\varphi_\alpha := \begin{cases} 1, & \text{falls } \bar{X}_n > \tilde{c}_\alpha, \\ 0, & \text{falls } \bar{X}_n < \tilde{c}_\alpha, \\ \gamma_\alpha, & \text{falls } \bar{X}_n = \tilde{c}_\alpha \end{cases}$$

mit  $\tilde{c}_\alpha \geq 0$ ,  $\gamma_\alpha \in [0, 1]$  geeignet. Unter  $\mathbb{P}_0$  und  $\mathbb{P}_1$  hat das Ereignis  $\{\bar{X}_n = \tilde{c}_\alpha\}$  die Wahrscheinlichkeit Null, sodass wir auf Randomisierung verzichten und einfach  $\gamma_\alpha = 0$  setzen können. Da unter  $\mathbb{P}_0$  für das Stichprobenmittel  $\bar{X}_n \sim N(\mu_0, \sigma^2/n)$  gilt, erhalten wir somit einen Neyman-Pearson-Test vom Niveau  $\alpha \in (0, 1)$  durch

$$\varphi_\alpha = \mathbb{1}(\bar{X}_n > \mu_0 + \sigma n^{-1/2} q_{1-\alpha})$$

mit dem  $(1 - \alpha)$ -Quantil  $q_{1-\alpha}$  der Standardnormalverteilung, die gerade mit dem einseitigen Gauß-Test aus Beispiel 1.56 übereinstimmt. Genau wie beim einseitigen

Binomialtest sieht man, dass  $\varphi_\alpha$  sogar ein Test vom Niveau  $\alpha$  für  $H_0: \mu \leq \mu_0$  gegen  $H_1: \mu > \mu_0$  mit UMP-Eigenschaft ist.

Wie diese Beispiele demonstrieren, lässt sich die Neyman-Pearson-Theorie manchmal auch vom Fall einfacher Hypothesen auf einseitige Testprobleme übertragen. Die notwendige Strukturvoraussetzung dafür waren immer monotone Likelihood-Quotienten. Für die wichtigen zweiseitigen Testprobleme führt dieser Ansatz jedoch nicht zum Ziel, man kann sogar zeigen, dass für zweiseitige Binomial- oder Gauß-Testprobleme keine UMP-Tests existieren können. Stattdessen sollte für zweiseitige Testprobleme die Klasse der unverfälschten Tests betrachtet werden, für die unter geeigneten Bedingungen die Existenz von *besten, unverzerrten Tests* (englisch: *uniformly most powerfull unbiased*, kurz: UMPU) gezeigt werden kann.

Auch für allgemeine zusammengesetzte Hypothesen führt uns die Neyman-Pearson-Theorie auf einen intuitiven Ansatz zur Wahl der Teststatistik. Nehmen wir dazu ein statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit Likelihood-Funktion  $L$  und Partition  $\Theta = \Theta_0 \cup \Theta_1$  an, so können wir eine Teststatistik für  $H_0: \vartheta \in \Theta_0$  gegen  $H_1: \vartheta \in \Theta_1$  konstruieren, indem wir in  $\Theta_0$  und  $\Theta_1$  jeweils den Parameter wählen, der die Likelihood-Funktion maximiert, und damit, dem Neymann-Pearson-Ansatz folgend, den entsprechenden verallgemeinerten Likelihood-Quotienten bilden:

$$T(x) := \frac{\sup_{\vartheta \in \Theta_1} L(\vartheta, x)}{\sup_{\vartheta \in \Theta_0} L(\vartheta, x)}, \quad (1.17)$$

wobei formal  $a/0 := +\infty$  gesetzt werde. Damit erhalten wir eine sehr allgemeine Methode, Tests zu konstruieren.

**Methode 1.63 (Likelihood-Quotiententest)** Für ein dominiertes statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  mit Likelihood-Funktion  $L(\vartheta, x)$  betrachte das Testproblem  $H_0: \vartheta \in \Theta_0$  gegen  $H_1: \vartheta \in \Theta_1$ . Für  $c_\alpha \geq 0$  und  $\gamma_\alpha \in [0, 1]$  sowie  $T$  aus (1.17) ist ein **Likelihood-Quotiententest** (englisch: *likelihood ratio test*, kurz: LR-Test) gegeben durch

$$\varphi_\alpha(x) = \mathbb{1}(T(x) > c_\alpha) + \gamma_\alpha \mathbb{1}(T(x) = c_\alpha), \quad x \in \mathcal{X}.$$

Es ist leicht einzusehen, dass unter  $H_i$ ,  $i = 0, 1$ , fast sicher  $\sup_{\vartheta \in \Theta_i} L(\vartheta) > 0$  gilt, sodass ein Likelihood-Quotiententest  $\varphi_\alpha$  fast sicher wohldefiniert ist (ist der Nenner in  $T$  gleich null, so setze  $T = +\infty$ , da dies fast sicher nur unter  $H_1$  geschieht, wo der Zähler fast sicher positiv ist).

Eine andere Interpretation des Likelihood-Quotiententests ergibt sich, wenn die Maximum-Likelihood-Schätzer  $\hat{\vartheta}_0$  und  $\hat{\vartheta}_1$  von  $\vartheta$  über den Parametermengen  $\Theta_0$  bzw.  $\Theta_1$  existieren. Dann ist der Likelihood-Quotient gerade

$$T(x) = \frac{L(\hat{\vartheta}_1(x), x)}{L(\hat{\vartheta}_0(x), x)},$$

und wir können den Likelihood-Quotiententest als einen Neyman-Pearson-Test zwischen den geschätzten Parametern  $\hat{\vartheta}_0$  und  $\hat{\vartheta}_1$  interpretieren. Beachte dazu aber,

dass die Schätzer zufällig sind und das Neyman-Pearson-Lemma keine Anwendung findet. Ähnlich wie bei der Maximum-Likelihood-Methode führt der Likelihood-Quotiententest häufig, aber nicht immer zu guten Tests.

*Beispiel 1.64 (Zweiseitiger Binomialtest)* In Beispiel 1.48 haben wir den zweiseitigen Binomialtest für  $\mathbb{P}_\vartheta = \text{Bin}(n, \vartheta)$  und  $H_0: \vartheta = \vartheta_0$  gegen  $H_1: \vartheta \neq \vartheta_0$  kennengelernt. Als Likelihood-Quotientenstatistik erhalten wir für  $\vartheta \in (0, 1)$  fest und  $\widehat{\vartheta}(k) = k/n$

$$T(k) = \frac{\sup_{\vartheta \neq \vartheta_0} \vartheta^k (1 - \vartheta)^{n-k}}{\vartheta_0^k (1 - \vartheta_0)^{n-k}} = \frac{\widehat{\vartheta}(k)^k (1 - \widehat{\vartheta}(k))^{n-k}}{\vartheta_0^k (1 - \vartheta_0)^{n-k}}, \quad k \in \{0, \dots, n\}.$$

Im Fall  $\vartheta_0 = 1/2$  vereinfacht sich dies durch Einsetzen zu  $T(k) = (2n)^{-n} k^k (n-k)^{n-k}$ . Damit ist  $T$  um  $n/2$  symmetrisch ( $T(k) = T(n-k)$  für  $0 \leq k \leq n/2$ ) und wachsend ( $T(k+1) > T(k)$  für  $k \geq n/2$ ). Der Likelihood-Quotiententest lässt sich somit schreiben als  $\varphi_\alpha(k) = \mathbb{1}(|k - n/2| > \widehat{c}_\alpha) + \gamma_\alpha \mathbb{1}(|k - n/2| = \widehat{c}_\alpha)$ , was genau dem zweiseitigen Binomialtest entspricht. Beachte, dass für  $\vartheta_0 \neq 1/2$  die Asymmetrie der  $\text{Bin}(n, \vartheta_0)$ -Verteilung zu einer Teststatistik führt, die nicht symmetrisch um  $\vartheta_0 n$  ist.

*Beispiel 1.65 (Zweiseitiger Gauß-Test)* Aufgrund einer mathematischen Stichprobe  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  mit bekanntem  $\sigma > 0$  wollen wir für festes  $\mu_0 \in \mathbb{R}$  die Hypothese  $H_0: \mu = \mu_0$  gegen  $H_1: \mu \neq \mu_0$  testen. Aus Stetigkeitsgründen ist das Supremum der Likelihood-Funktion über  $\mu \in \mathbb{R} \setminus \{\mu_0\}$  gleich dem Supremum über ganz  $\mathbb{R}$ , was am Stichprobenmittel als Maximum-Likelihood-Schätzer  $\widehat{\mu} = \bar{X}$  angenommen wird, und die Likelihood-Quotientenstatistik (in den Beobachtungen geschrieben) ist

$$T = \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n \left( (X_i - \bar{X})^2 - (X_i - \mu_0)^2 \right) \right) = \exp \left( \frac{n}{2\sigma^2} (\mu_0 - \bar{X})^2 \right).$$

Eine einfache statistische Herleitung der zweiten Identität ergibt sich, wenn man

$$\frac{1}{n} \sum_{i=1}^n (X_i - \mu_0)^2 = (\bar{X} - \mu_0)^2 + \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

als Bias-Varianz-Zerlegung bezüglich der empirischen Verteilung der  $(X_i)$  versteht. Wir bemerken, dass  $T$  in  $|\bar{X} - \mu_0|$  streng monoton wächst und auf Randomisierung verzichtet werden kann, sodass

$$\varphi_\alpha = \mathbb{1}(|\bar{X} - \mu_0| > \sigma n^{-1/2} q_{1-\alpha/2})$$

wegen  $\bar{X} \sim N(\mu_0, \sigma^2/n)$  unter  $H_0$  ein Likelihood-Quotiententest zum Niveau  $\alpha$  ist. Man nennt  $\varphi_\alpha$  *zweiseitigen Gauß-Test*.

## 1.4 Konfidenzmengen

In unserem Saskia-Beispiel 1.7 wird eine Umfrage unter Studierenden durchgeführt, die nur von einem Teil der Studierenden beantwortet wird. Die Antworten aus der Stichprobe ergeben im Mittel einen Wert, von dem wir nicht wissen, ob er der Durchschnittsantwort aller Studierenden (auch jener, die nicht an der Umfrage teilgenommen haben) entspricht. Anhand dieser Stichprobe lässt sich jedoch ein Intervall angeben, das den wahren Mittelwert mit einer gegebenen Wahrscheinlichkeit enthält, und zwar egal, welcher wahre Mittelwert die Wahrscheinlichkeitsverteilung bestimmt.

Statt wie in der Parameterschätzung einen einzelnen Wert zu bestimmen, der möglichst in der Nähe des wahren Parameters liegt, wollen wir also nun Bereiche angeben, von denen wir mit einer gewissen Zuversicht (Konfidenz) sagen können, dass der wahre Parameter in ihnen liegt. Daher werden diese Bereiche auch als „Konfidenzmengen“ oder „Konfidenzbereiche“ bezeichnet. In allen Wissenschaften ist diese Quantifizierung der statistischen Unsicherheit (*uncertainty quantification*) von großer Bedeutung.

**Definition 1.66** Sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell mit abgeleitetem Parameter  $\rho: \Theta \rightarrow \mathbb{R}^d$ . Eine mengenwertige Abbildung

$$C: \mathcal{X} \rightarrow \mathcal{P}(\mathbb{R}^d)$$

heißt **Konfidenzmenge zum Konfidenzniveau**  $1 - \alpha$  (oder zum Irrtumsniveau  $\alpha$ ) für  $\alpha \in (0, 1)$ , falls die Messbarkeitsbedingung  $\{x \in \mathcal{X} : \rho(\vartheta) \in C(x)\} \in \mathcal{F}$  für alle  $\vartheta \in \Theta$  erfüllt ist und

$$\mathbb{P}_\vartheta(\rho(\vartheta) \in C) = \mathbb{P}_\vartheta(\{x \in \mathcal{X} : \rho(\vartheta) \in C(x)\}) \geq 1 - \alpha \quad \text{für alle } \vartheta \in \Theta$$

gilt. Im Fall  $d = 1$  und falls  $C(x)$  für jedes  $x \in \mathcal{X}$  ein Intervall ist, heißt  $C$  **Konfidenzintervall**.

*Bemerkung 1.67 (Konfidenzmengen)* Hier ist  $\rho(\vartheta)$  fixiert, während  $C$  zufällig ist. Konfidenzmengen sind also wie folgt zu interpretieren: Werden in  $m$  unabhängigen Experimenten für (verschiedene) Parameter Konfidenzmengen zum Niveau  $1 - \alpha = 0,95$  konstruiert, dann liegt der unbekannte Parameter  $\rho(\vartheta)$  in 95% der Fälle in der jeweiligen Konfidenzmenge (für  $m$  hinreichend groß; starkes Gesetz der großen Zahlen), unabhängig davon, welches  $\vartheta \in \Theta$  vorliegt. Man spricht aber *nicht* davon, dass mit 95%-iger Wahrscheinlichkeit  $\rho(\vartheta)$  in  $C(x)$  liegt, denn wir betrachten gar kein Wahrscheinlichkeitsmaß auf  $\Theta$ .

Die Menge  $C = \mathbb{R}^d$  ist stets eine (triviale) Konfidenzmenge zu jedem beliebigen Niveau, die uns aber keine Information über  $\rho(\vartheta)$  liefert. Je kleiner die Konfidenzmenge ist, desto genauere Aussagen erhalten wir über den unbekannten Parameter. Wir werden uns daher bemühen, Konfidenzmengen möglichst klein zu konstruieren, sodass das Konfidenzniveau gerade noch eingehalten wird.

Wie kann von einer erhobenen Stichprobe ausgehend eine Konfidenzmenge konstruiert werden? Ein häufig verwendetes Konstruktionsprinzip für die Konfidenzintervalle ist die Verwendung eines Schätzers und seiner Verteilung, wie die nächsten Beispiele illustrieren.

*Beispiel 1.68 (Konstruktion von Konfidenzintervallen – a)* Wir wollen ein Konfidenzintervall für Saskias Umfrage zur Studierendenzufriedenheit konstruieren. In Beispiel 1.7 hatten wir das statistische Modell  $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\text{Ber}(p))^{\otimes n})_{p \in (0, 1)}$  betrachtet und als Schätzer  $\hat{\rho}_n$  wählten wir das arithmetische Mittel der Stichprobe. Unsere Grundidee zur Konstruktion des Konfidenzintervalls ist es, um den Schätzer  $\hat{\rho}_n$  ein symmetrisches Intervall

$$C_n := [\hat{\rho}_n - \varepsilon_n, \hat{\rho}_n + \varepsilon_n]$$

aufzuspannen, wobei wir  $\varepsilon_n$  noch näher bestimmen müssen. Damit  $C_n$  ein Konfidenzintervall zum Irrtumsniveau  $\alpha \in (0, 1)$  ist, fordern wir

$$\mathbb{P}_p(p \in C_n) = \mathbb{P}_p(|\hat{\rho}_n - p| \leq \varepsilon_n) = \mathbb{P}_p\left(\left|\sum_{i=1}^n (X_i - p)\right| \leq n\varepsilon_n\right) \stackrel{!}{\geq} 1 - \alpha.$$

Wegen  $\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$  können wir  $\varepsilon_n$  numerisch mithilfe der Quantile der Binomialverteilung bestimmen. Alternativ dazu wird im Folgenden eine Normalapproximation verwendet, wobei das resultierende Konfidenzintervall dann nur asymptotisch für große  $n$  das Niveau  $1 - \alpha$  besitzt.

Es gilt für eine Zufallsvariable  $Z \sim N(0, 1)$

$$\begin{aligned} \mathbb{P}_p(p \in C_n) &= \mathbb{P}_p\left(\frac{|\sum_i X_i - np|}{\sqrt{np(1-p)}} \leq \frac{n\varepsilon_n}{\sqrt{np(1-p)}}\right) \\ &\approx \mathbb{P}\left(|Z| \leq \sqrt{\frac{n}{p(1-p)}} \varepsilon_n\right) \\ &= \Phi\left(\sqrt{\frac{n}{p(1-p)}} \varepsilon_n\right) - \Phi\left(-\sqrt{\frac{n}{p(1-p)}} \varepsilon_n\right) \\ &= 2\Phi\left(\sqrt{\frac{n}{p(1-p)}} \varepsilon_n\right) - 1 \stackrel{!}{=} 1 - \alpha, \end{aligned}$$

wobei wir die Gleichheit mit  $1 - \alpha$  fordern, um das Konfidenzniveau voll auszuschöpfen. Mit dem  $(1 - \alpha/2)$ -Quantil  $q_{1-\alpha/2}$  erhalten wir also mit gegen  $1 - \alpha$  konvergierender Wahrscheinlichkeit

$$|\hat{\rho}_n - p| \leq \sqrt{\frac{p(1-p)}{n}} q_{1-\alpha/2}. \quad (1.18)$$

Die rechte Seite ist jedoch noch nicht zur Konstruktion eines Konfidenzintervalls geeignet, da sie vom unbekannten  $p$  abhängt. Wir ersetzen  $p$  an dieser Stelle daher durch den Schätzer  $\hat{\rho}_n$  und wählen  $\varepsilon_n = \sqrt{\frac{\hat{\rho}_n(1-\hat{\rho}_n)}{n}} q_{1-\alpha/2}$ . Damit ergibt sich das

sogenannte *Wald-Intervall*

$$C_n = \left[ \hat{\rho}_n - \sqrt{\frac{\hat{\rho}_n(1-\hat{\rho}_n)}{n}} q_{1-\alpha/2}, \hat{\rho}_n + \sqrt{\frac{\hat{\rho}_n(1-\hat{\rho}_n)}{n}} q_{1-\alpha/2} \right].$$

In der Tat gilt für jedes  $\tau > 0$

$$\begin{aligned} \mathbb{P}_p(p \notin C_n) &= \mathbb{P}_p\left(|\hat{\rho}_n - p| > \sqrt{\frac{\hat{\rho}_n(1-\hat{\rho}_n)}{n}} q_{1-\alpha/2}\right) \\ &\leq \mathbb{P}_p\left(|\hat{\rho}_n - p| > \sqrt{\frac{p(1-p)}{n}} q_{1-\alpha/2}(1-\tau)\right) \\ &\quad + \mathbb{P}_p\left(\sqrt{p(1-p)} - \sqrt{\hat{\rho}_n(1-\hat{\rho}_n)} > \tau\sqrt{p(1-p)}\right). \end{aligned}$$

Für  $n \rightarrow \infty$  konvergiert der erste Term wie in obiger Rechnung gegen  $2 - 2\Phi(q_{1-\alpha/2}(1-\tau))$ , während der zweite Term aufgrund der Konsistenz  $\hat{\rho}_n \xrightarrow{\mathbb{P}_p} p$  (Gesetz der großen Zahlen) für jedes  $\tau$  gegen 0 konvergiert. Aus  $\lim_{\tau \rightarrow 0} \Phi(q_{1-\alpha/2}(1-\tau)) = 1 - \alpha/2$  ergibt sich damit  $\lim_{n \rightarrow \infty} \mathbb{P}_p(p \notin C_n) = \alpha$ , sodass  $C_n$  tatsächlich ein Konfidenzintervall zum asymptotischen Niveau  $1 - \alpha$  ist.

Alternativ kann man auch (1.18) durch Quadrieren und Auflösen der quadratischen Gleichung direkt nach  $p$  umstellen. Das vermeidet den zusätzlichen Approximationsschritt und führt auf das *Wilson-Intervall* (Aufgabe 1.13).

*Beispiel 1.69 (Konstruktion von Konfidenzintervallen – b)* Selbst wenn die Verteilung des Schätzers  $\hat{\rho}_n$  nicht explizit bekannt ist, können Konfidenzintervalle konstruiert werden, sofern sich  $\hat{\rho}_n$  um den wahren Wert  $\rho(\vartheta)$  konzentriert. In diesem Fall verwenden wir, dass die Wahrscheinlichkeit  $\mathbb{P}_\vartheta(|\hat{\rho}_n - \rho(\vartheta)| > \varepsilon)$  klein wird, wenn  $\varepsilon$  groß wird. Wir wollen diese Idee im Modell aus Beispiel 1.68 verdeutlichen. Für  $C_n := (\hat{\rho}_n - \varepsilon_n, \hat{\rho}_n + \varepsilon_n)$  gilt wegen der Tschebyscheff-Ungleichung

$$\begin{aligned} \mathbb{P}_p(p \in C_n) &= \mathbb{P}_p(|\hat{\rho}_n - p| < \varepsilon_n) \\ &= \mathbb{P}_p\left(\left|\sum_{i=1}^n X_i - np\right| < n\varepsilon_n\right) \\ &\geq 1 - \frac{\text{Var}(\sum_i X_i)}{n^2\varepsilon_n^2} \\ &= 1 - \frac{np(1-p)}{n^2\varepsilon_n^2} \stackrel{!}{\geq} 1 - \alpha. \end{aligned}$$

Wir schätzen  $p(1-p) \leq 1/4$  ab und erhalten  $\varepsilon_n = 1/\sqrt{4n\alpha}$ . Damit erhalten wir das (nicht asymptotische) Konfidenzintervall  $C_n = \left(\hat{\rho}_n - \frac{1}{2\sqrt{n\alpha}}, \hat{\rho}_n + \frac{1}{2\sqrt{n\alpha}}\right)$ . Allerdings ist die Abschätzung mit der Tschebyscheff-Ungleichung sehr grob und führt daher zu einem sehr vorsichtigen bzw. konservativen Konfidenzintervall.

Eine alternative Konstruktion von Konfidenzmengen bietet folgender Korrespondenzsatz:

**Satz 1.70 (Korrespondenzsatz)** Es sei  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $\alpha \in (0, 1)$ . Dann gilt:

- (i) Liegt für jedes  $\vartheta_0 \in \Theta$  ein nicht-randomisierter Test  $\varphi_{\vartheta_0}$  der Hypothese  $H_0: \vartheta = \vartheta_0$  zum Signifikanzniveau  $\alpha$  vor, so definiert  $C(x) := \{\vartheta_0 \in \Theta : \varphi_{\vartheta_0}(x) = 0\}$  für  $x \in X$  eine Konfidenzmenge für  $\vartheta$  zum Konfidenzniveau  $1 - \alpha$ .
- (ii) Ist  $C$  eine Konfidenzmenge für  $\vartheta$  zum Niveau  $1 - \alpha$ , dann ist  $\varphi_{\vartheta_0}(x) := 1 - \mathbb{1}_{C(x)}(\vartheta_0)$  ein Niveau- $\alpha$ -Test der Hypothese  $H_0: \vartheta = \vartheta_0$ .

Es gibt also eine Eins-zu-eins-Beziehung zwischen Hypothesentests und Konfidenzmengen. Die Konfidenzmenge in (i) enthält all jene  $\vartheta_0$ , für die der Test  $\varphi_{\vartheta_0}$  die Nullhypothese aufgrund der beobachteten Stichprobe  $x$  nicht verwirft.

**Beweis** Nach Konstruktion erhält man in beiden Fällen

$$\forall \vartheta \in \Theta : \forall x \in X : \varphi_\vartheta(x) = 0 \iff \vartheta \in C(x).$$

Damit ist  $\varphi_\vartheta$  genau dann ein Test zum Niveau  $\alpha$  für alle  $\vartheta$ , wenn

$$1 - \alpha \leq \mathbb{P}_\vartheta(\varphi = 0) = \mathbb{P}_\vartheta(\{x : \vartheta \in C(x)\}),$$

und somit ist  $C$  eine Konfidenzmenge zum Niveau  $\alpha$ . □

Bezeichnen wir den Annahmebereich der Tests  $\varphi_\vartheta$  mit  $A(\vartheta)$ , liefert uns der Korrespondenzsatz folgende Methode:

**Methode 1.71 (Konstruktion von Konfidenzmengen)** Sei  $(X, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein statistisches Modell und  $\alpha \in (0, 1)$ . Wähle zu jedem  $\vartheta \in \Theta$  ein  $A(\vartheta) \in \mathcal{F}$  mit  $\mathbb{P}_\vartheta(A(\vartheta)) \geq 1 - \alpha$  und setze  $C(x) := \{\vartheta \in \Theta : x \in A(\vartheta)\}$  für  $x \in X$ . Dann gilt

$$x \in A(\vartheta) \iff \vartheta \in C(x),$$

woraus mit obigem Satz folgt

$$\mathbb{P}_\vartheta(\{x \in X : \vartheta \in C(x)\}) \geq 1 - \alpha \quad \forall \vartheta \in \Theta.$$

Man könnte zunächst meinen, dass durch diese Methode die Schwierigkeit der Konstruktion lediglich von  $C(x)$  auf  $A(\vartheta)$  verschoben wurde, aber dadurch haben wir einen Vorteil erlangt:  $A(\vartheta)$  ist eine Teilmenge von  $X$ , die wir mit  $\mathbb{P}_\vartheta$  messen können, und aus der Kenntnis von  $\mathbb{P}_\vartheta$  ergibt sich meist eine einfache Wahl eines Ereignisses  $A(\vartheta)$  mit Wahrscheinlichkeit  $1 - \alpha$ .

**Beispiel 1.72 (Konstruktion von Konfidenzmengen)** Wir wollen ein Konfidenzintervall für die Geburtswahrscheinlichkeit von Mädchen in Hamburg berechnen. In Beispiel 1.48 hatten wir die Anzahl der Mädchen unter  $n$  Geburten mit  $X \sim \text{Bin}(n, \vartheta)$  modelliert, wobei  $\vartheta \in \Theta := (0, 1)$  und der Stichprobenraum  $X = \{0, 1, \dots, n\}$  ist. Wir wählen ein Niveau  $\alpha \in (0, 1)$ . Die Annahmebereiche der in Beispiel 1.48 konstruierten Tests sind gegeben durch

$$A(\vartheta) := \{x \in X : u_\alpha(\vartheta) \leq x \leq o_\alpha(\vartheta)\}$$

mit  $u_\alpha(\vartheta)$  und  $o_\alpha(\vartheta)$  aus (1.11). Damit schneiden wir von der gesamten möglichen Wahrscheinlichkeitsmasse an beiden Enden  $\frac{\alpha}{2}$  ab, sodass in der Mitte  $1 - \alpha$  übrig bleibt (siehe Abbildung 1.2). Da wir  $A(\vartheta)$  gewählt haben, nutzen wir nun die Beziehung

$$x \in A(\vartheta) \iff \vartheta \in C(x) \quad \text{für alle } x \in \mathcal{X}, \vartheta \in \Theta,$$

um die Konfidenzmenge  $C$  zu bestimmen. Die Leserin überprüfe selbst, dass  $u_\alpha: [0, 1] \rightarrow \{0, 1, \dots, n\}$  monoton fallend und rechtsseitig stetig ist und  $o_\alpha: [0, 1] \rightarrow \{0, 1, \dots, n\}$  monoton wachsend und linksseitig stetig ist. Folglich ist

$$C(x) := [\inf\{\vartheta \in \Theta : o_\alpha(\vartheta) = x\}, \sup\{\vartheta \in \Theta : u_\alpha(\vartheta) = x\}]$$

ein Konfidenzintervall zum Konfidenzniveau  $1 - \alpha$ . Es wird *Clopper-Pearson-Intervall* genannt. Die Berechnung des Infimums und des Supremums ist eine numerische Aufgabe.

Für  $n = 21.126$  Geburten im Jahr 2018, von denen 10.215 weiblich waren, erhalten wir folgendes Konfidenzintervall zum Niveau  $1 - \alpha = 0,95$  für die Wahrscheinlichkeit, dass ein Mädchen geboren wurde:

$$C = [0,4768; 0,4903]$$

**Bemerkung 1.73 (Einseitige Konfidenzbereiche)** Bisher haben wir nur zweiseitige Konfidenzintervalle gesehen. Eine andere Variante sind *einseitige Konfidenzintervalle*. Das heißt, dass nur eine Seite von den Beobachtungen abhängt und die andere fest ist. Im vorherigen Beispiel könnte man analog das Konfidenzintervall

$$\tilde{C}(x) := [0, \sup\{\vartheta \in \Theta : \tilde{u}_\alpha(\vartheta) = x\}), \quad x \in \mathcal{X},$$

konstruieren mit

$$\tilde{u}_\alpha(\vartheta) := \max\{k \in \mathcal{X} : \mathbb{P}_\vartheta(X < k) \leq \alpha\}.$$

Nach unten verliert diese Konstruktion des Konfidenzintervalls zwar an Aussagekraft, aber nach oben gewinnt es ebenjene, da die Obergrenze schärfer wird.

## 1.5 Aufgaben

1.1 Wir betrachten eine auf dem Intervall  $[a, b]$  gleichverteilte mathematische Stichprobe  $X_1, \dots, X_n$  für  $n \in \mathbb{N}$  und unbekannte Parameter  $-\infty < a < b < \infty$ .

- Formalisieren Sie das statistische Modell.
- Bestimmen Sie Momentenschätzer für  $a$  und  $b$ .
- Bestimmen Sie die Maximum-Likelihood-Schätzer für  $a$  und  $b$ .
- Welches quadratische Risiko hat der Maximum-Likelihood-Schätzer?

- 1.2 Um die Gesamtanzahl  $N$  der insgesamt in Berlin registrierten Taxis zu schätzen, notiert sich ein Tourist die Konzessionsnummern von  $n < N$  vorbeifahrenden Taxis (Wiederholungen möglich). Er nimmt an, dass alle Taxis von 1 bis  $N$  durchnummeriert sind und die beobachteten Taxinummern unabhängig voneinander und mit gleicher Wahrscheinlichkeit vorbeifahren.
- Formalisieren Sie das statistische Modell.
  - Berechnen Sie aus den notierten Nummern  $X_1, \dots, X_n$  einen Maximum-Likelihood-Schätzer  $\hat{N}$  für  $N$ . Ist dieser erwartungstreu?
  - Berechnen Sie approximativ für großes  $N$  den relativen Erwartungswert  $\mathbb{E}[\hat{N}]/N$ .
- Hinweis:* Fassen Sie einen geeigneten Ausdruck als Riemann-Summe auf.
- 1.3 Es sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  ein dominiertes statistisches Modell mit zwei dominierenden Maßen  $\mu_1$  und  $\mu_2$ . Die zugehörigen Likelihoodfunktionen seien mit  $L_1(\vartheta, x) = \frac{d\mathbb{P}_\vartheta}{d\mu_1}(x)$  bzw.  $L_2(\vartheta, x) = \frac{d\mathbb{P}_\vartheta}{d\mu_2}(x)$ ,  $x \in \mathcal{X}$ ,  $\vartheta \in \Theta$ . Zeigen Sie, dass jeder Maximum-Likelihood-Schätzer unabhängig vom dominierenden Maß ist, d.h. jede Maximalstelle von  $L_1$  ist auch eine Maximalstelle von  $L_2$ .
- 1.4 Sei  $X_1, \dots, X_n$  eine mathematische Stichprobe reeller Zufallsvariablen mit der Verteilung  $\mathbb{P}_\vartheta$ ,  $\vartheta \in \Theta$ . Für  $k > 0$  betrachten wir die Funktion

$$\psi_k: \mathbb{R} \rightarrow \mathbb{R}, \quad \psi_k(x) = \begin{cases} -k, & x < -k, \\ x, & |x| \leq k, \\ k, & x > k. \end{cases}$$

- Beweisen Sie, dass ein  $h_k \in \mathbb{R}$  existiert, sodass  $\mathbb{E}[\psi_k(X_1 - h_k)] = 0$ . Der Huber-Schätzer  $\hat{h}_k$  von  $h_k$  ist definiert als Nullstelle der Funktion

$$\mathbb{R} \ni h \mapsto \sum_{i=1}^n \psi_k(X_i - h).$$

Weisen Sie nach, dass auch  $\hat{h}_k$  stets existiert.

- Zeigen Sie, dass  $h_k$  für  $k \rightarrow 0$  gegen einen Median von  $\mathbb{P}_\vartheta$  konvergiert. Beweisen Sie im Fall  $\mathbb{E}[|X_1|] < \infty$ , dass  $h_k$  für  $k \rightarrow \infty$  gegen den Mittelwert von  $\mathbb{P}_\vartheta$  konvergiert.
- Sei nun  $\mathbb{P}_\vartheta$  die Cauchy-Verteilung mit Parameter  $\vartheta = (x_0, \gamma) \in \mathbb{R} \times \mathbb{R}_+$  und Lebesgue-Dichte

$$f_{x_0, \gamma}(x) = \frac{1}{\pi} \frac{\gamma}{(x - x_0)^2 + \gamma^2}, \quad x \in \mathbb{R}.$$

Bestimmen Sie  $h_k$  in Abhängigkeit von  $\vartheta = (x_0, \gamma)$  für alle  $k > 0$ .

- Simulieren Sie  $n = 200$  unabhängige Zufallsvariablen mit den Randverteilungen

- (i) Lognormalverteilung  $\log N(\mu, \sigma)$  mit Parametern  $\mu = 1, \sigma^2 = 1$  sowie  $\mu = 1, \sigma^2 = 3$ .
- (ii) Cauchy-Verteilung mit Parametern  $x_0 = 0, \gamma = 1$  und  $x_0 = 0, \gamma = 5$ .  
Bestimmen Sie in allen vier Szenarien den theoretischen Median und den Mittelwert. Berechnen Sie außerdem  $\hat{h}_k$  für  $k \in \{\frac{1}{10}, \frac{1}{2}, 1, 2, 5, 10, 20, 50\}$ . Stellen Sie Ihre Ergebnisse graphisch dar.
- 1.5 Sei  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\theta)_{\theta \in \Theta})$  mit  $\Theta = \mathbb{R}$  ein von  $\mu$  dominiertes statistisches Modell mit Dichten  $f_{X|T=\theta} := \frac{d\mathbb{P}_\theta}{d\mu}$  und sei  $\pi$  eine a-priori-Verteilung auf  $(\Theta, \mathcal{F}_\Theta)$  mit Dichte  $f_T$  bezüglich eines Maßes  $\nu$ . Nehmen Sie an, dass  $f_{X|T=\cdot} : \mathcal{X} \times \Theta \rightarrow \mathbb{R}_+$   $(\mathcal{F} \otimes \mathcal{F}_\Theta)$ -messbar ist. Zeigen Sie:
- Unter quadratischem Verlust ist der Bayes-Schätzer gegeben durch den a-posteriori-Erwartungswert.
  - Unter absolutem Verlust ist der Bayes-Schätzer gegeben durch den a-posteriori-Median.
  - Unter 0-1-Verlust ist der Bayes-Schätzer gegeben durch den a-posteriori-Modus (das heißt der Maximalstelle der a-posteriori-Verteilung).
- 1.6 Es sei  $X \sim \text{Bin}(n, p)$  binomialverteilt mit  $n \in \mathbb{N}$  und  $p \in [0, 1]$ , und sei die a-priori-Verteilung  $\pi$  für  $p$  auf  $[0, 1]$  gegeben durch eine Beta-Verteilung  $\text{Beta}(\alpha, \beta)$  mit Parametern  $\alpha, \beta > 0$ .
- Zeigen Sie, dass die Gleichverteilung  $U([0, 1])$  ein Spezialfall der Beta-Verteilung ist.
  - Beweisen Sie, dass die Beta-Verteilungen zur Binomialverteilung konjugiert sind, das heißt die a-posteriori-Verteilung ist wieder Beta-verteilt. Bestimmen Sie die Parameter der a-posteriori-Beta-Verteilung.
  - Folgern Sie, dass der Bayes-Schätzer unter quadratischem Verlust gegeben ist durch
 
$$\hat{p}_{a,b} = \frac{a + X}{a + b + n}.$$
 Bestimmen Sie dessen mittleres quadratisches Risiko  $\mathbb{E}_p[|\hat{p}_{a,b} - p|^2]$  in Abhängigkeit von  $a, b$  und  $p$ .
  - Wählen Sie  $a^*, b^* > 0$  so, dass  $\max_{p \in [0,1]} \mathbb{E}_p[|\hat{p}_{a,b} - p|^2]$  minimal ist. Folgern Sie, dass  $\hat{p}_{a^*,b^*}$  minimax-Schätzer ist und der Maximum-Likelihood-Schätzer  $\hat{p} = X/n$  nicht minimax ist.
- 1.7 Die Vertriebsleiterin einer Getreidemühle geht davon aus, dass ihre angebotenen Mehlpackungen ein mittleres Füllgewicht von 1000 g mit einer Standardabweichung von 12,5 g haben. Sie möchte die Füllmenge mit einem statistischen Test überprüfen. Eine genauere Untersuchung an  $n = 40$  Packungen zeigt, dass die Füllmenge im Durchschnitt 1004,32 g beträgt.
- Formalisieren Sie das statistische Modell und das Testproblem unter einer Normalverteilungsannahme an die Packungsgewichte.

- (b) Konstruieren Sie einen Test  $\varphi$  zum Niveau  $\alpha = 0,05$  unter Verwendung des Durchschnittsgewichts als Teststatistik. Verwirft dieser Test die Hypothese?
- (c) Stellen Sie die Gütefunktion von  $\varphi$  mithilfe der Verteilungsfunktion der Standardnormalverteilung  $\Phi: \mathbb{R} \rightarrow [0, 1]$  dar.

- 1.8 Am Hamburger U-Bahnhof Schlump treffen  $n \in \mathbb{N}$  Studierende ein und müssen jeweils  $X_i$  Minuten auf die Bahn warten,  $i = 1, \dots, n$ . Bezeichnet  $\vartheta > 0$  die erwartete Wartezeit, soll

$$H_0: \vartheta \leq \vartheta_0 \quad \text{gegen} \quad H_1: \vartheta > \vartheta_0$$

für ein  $\vartheta_0 > 0$  getestet werden. Hierzu wird die Teststatistik  $\varphi_c(X_1, \dots, X_n) = \mathbb{1}_{(c, \infty)}(\bar{X}_n)$  mit kritischem Wert  $c > 0$  und  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  verwendet.

- (a) Begründen Sie, warum die Exponentialverteilungsannahme  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Exp}(\lambda)$  mit Parameter  $\lambda > 0$  sinnvoll ist, und formulieren Sie das Testproblem in Abhängigkeit von  $\lambda$ .
  - (b) Bestimmen Sie die Verteilung von  $\bar{X}_n$ . *Hinweis:*  $\text{Exp}(\lambda) = \Gamma(1, \lambda)$ .
  - (c) Berechnen Sie  $c > 0$  so, dass  $\varphi_c$  ein Test zum Niveau  $\alpha \in (0, 1)$  ist.
  - (d) Notieren Sie sich die Zeit, die Sie heute auf dem Heimweg auf die U-Bahn warten mussten.
- 1.9 Betrachten Sie ein statistisches Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$  und einen Test der Hypothese  $H_0: \vartheta \in \Theta_0 \neq \emptyset$  zum Niveau  $\alpha \in (0, 1)$  der Form  $\varphi = \mathbb{1}_{\{T > c_\alpha\}}$  mit Teststatistik  $T$  und kritischen Werten  $c_\alpha, \alpha \in (0, 1)$ . Falls  $\sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T > c_\alpha) \leq \alpha$  für alle  $\alpha \in (0, 1)$ , gilt für den p-Wert

$$p_\varphi(x) = \inf \{ \alpha \in [0, 1] : T(x) > c_\alpha \}.$$

- 1.10 Zeigen Sie, dass im statistischen Modell  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \{0,1\}})$  zu jedem  $\alpha \in (0, 1)$  ein Neyman-Pearson-Test zum Niveau  $\alpha$  existiert.
- 1.11 Die mathematische Stichprobe  $X_1, \dots, X_n$  sei gemäß  $N(\mu, \sigma^2)$  verteilt mit unbekanntem  $\mu \in \mathbb{R}$  und unbekanntem  $\sigma > 0$ .
- (a) Bestimmen Sie den Maximum-Likelihood-Schätzer  $\hat{\sigma}^2$  für  $\sigma^2$ .
  - (b) Konstruieren Sie einen Likelihood-Quotiententest für das Testproblem

$$H_0: \sigma^2 = \sigma_0^2 \quad \text{gegen} \quad H_1: \sigma^2 \neq \sigma_0^2$$

für ein  $\sigma_0 > 0$ . Zeigen Sie, dass dieser Test für geeignete  $a, b > 0$  von der Form  $T = 1 - \mathbb{1}(a \leq \frac{\hat{\sigma}^2}{\sigma_0^2} \leq b)$  ist.

- (c) Verwenden Sie, dass  $\frac{n\hat{\sigma}^2}{\sigma_0^2}$  unter der Hypothese  $\chi^2$ -verteilt ist, um  $a$  und  $b$  so zu wählen, dass  $T$  ein Test zum Niveau  $\alpha \in (0, 1)$  ist.
- 1.12 Ein Experimentator macht  $n \in \mathbb{N}$  unabhängige normalverteilte Messungen mit unbekanntem Erwartungswert  $\mu \in \mathbb{R}$ . Die Varianz  $\sigma^2 > 0$  meint er zu kennen.

- (a) Welches Konfidenzintervall  $C_{\alpha, \sigma^2}$  für  $\mu$  wird er zu einem vorgegeben Irrtumsniveau  $\alpha \in (0, 1)$  angeben?
- (b) Welches Irrtumsniveau hat dieses Konfidenzintervall  $C_{\alpha, \sigma^2}$ , wenn die Varianz in Wirklichkeit den Wert  $\tilde{\sigma}^2 > 0$  annimmt?
- (c) Simulieren Sie in 500 Durchgängen jeweils 20 unabhängige  $N(0, 1)$ -verteilte Zufallsvariablen  $X_1, \dots, X_{20}$  bzw.  $N(0, 2)$ -verteilte Zufallsvariablen  $Y_1, \dots, Y_{20}$ . Bestimmen Sie in jeder Iteration  $C_{\alpha, 1}$  unter Verwendung von  $(X_i)$  bzw.  $(Y_i)$ . Für wie viele Realisierungen von  $(X_i)$  bzw.  $(Y_i)$  liegt  $\mu = 0$  in  $C_{\alpha, 1}$  für  $\alpha \in \{\frac{1}{100}, \frac{1}{20}, \frac{1}{10}, \frac{1}{4}, \frac{1}{2}\}$ . Bestimmen Sie jeweils das Verhältnis zur Anzahl der Simulationsdurchgänge 500.

- 1.13 Es seien  $X_1, \dots, X_n$  unabhängige  $\text{Ber}(p)$  verteilte Zufallsvariablen mit unbekanntem Erfolgsparameter  $p$  und dem Schätzer  $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Weisen Sie nach, dass für  $\alpha \in (0, 1)$  das *Wilson-Intervall*  $[u_\alpha, o_\alpha]$  mit

$$u_\alpha = \frac{1}{1 + \frac{q_{1-\alpha/2}^2}{n}} \left( \hat{p}_n + \frac{q_{1-\alpha/2}^2}{2n} - \frac{q_{1-\alpha/2}}{\sqrt{n}} \sqrt{\hat{p}_n(1 - \hat{p}_n) + \frac{q_{1-\alpha/2}^2}{4n}} \right),$$

$$o_\alpha = \frac{1}{1 + \frac{q_{1-\alpha/2}^2}{n}} \left( \hat{p}_n + \frac{q_{1-\alpha/2}^2}{2n} + \frac{q_{1-\alpha/2}}{\sqrt{n}} \sqrt{\hat{p}_n(1 - \hat{p}_n) + \frac{q_{1-\alpha/2}^2}{4n}} \right).$$

ein Konfidenzintervall zum asymptotischen Konfidenzniveau  $1 - \alpha$  ist.

## Kapitel 2

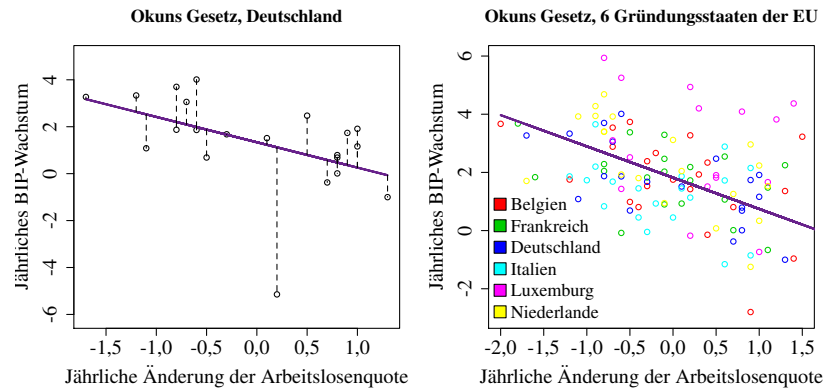
# Das lineare Modell

Mit der einfachen linearen Regression beginnend, werden wir in diesem Kapitel das lineare Modell im Detail studieren. Dieses schließt die multiple sowie die polynomiale Regression ein. Insbesondere befassen wir uns mit der verallgemeinerten Methode der kleinsten Quadrate. Unter einer Normalverteilungsannahme an die Beobachtungsfehler konstruieren wir anschließend Hypothesentests und Konfidenzmengen basierend auf t- und F-Statistiken. Als Spezialfall ergibt sich die Varianzanalyse.

### 2.1 Regression und kleinste Quadrate

Nachdem wir uns bisher mit einigen Grundbegriffen und Fragestellungen in der Statistik befasst haben, wollen wir nun das lineare Modell studieren. Die beobachteten Daten werden in dieser Modellklasse durch bekannte Einflussvariablen erklärt, wobei eine lineare Abhängigkeit von den unbekannten Parametern angenommen wird. Lineare Modelle finden unzählige Anwendungen und sind mathematisch relativ leicht zu analysieren. Insbesondere die Methode der kleinsten Quadrate, auf die wir noch genauer eingehen werden, hat seit über 200 Jahren nicht an Bedeutung verloren.

*Regression* bezeichnet die statistische Analyse des (nicht unbedingt linearen) Zusammenhangs zwischen einer *Zielgröße*  $Y$  (auch *Regressand*, *Response-Variable* oder *abhängige Variable* genannt) und einem Vektor von *Kovariablen*  $X = (X_1, \dots, X_n)$  (oder auch *Regressoren*, *erklärenden Variablen*, *unabhängigen Variablen*). Die Kovariablen können zufällig oder deterministisch sein. Wir sprechen in diesen Fällen von *zufälligem Design* bzw. *deterministischem Design*. Im Folgenden betrachten wir zunächst den letzteren Fall und schreiben der Klarheit halber  $X = (x_1, \dots, x_n)$ . Sind die Kovariablen zufällig, aber von den Beobachtungsfehlern unabhängig, so können wir auf  $X$  bedingen und die Resultate von deterministischem Design auf zufälliges Design übertragen.



**Abb. 2.1** Jährliche prozentuale Veränderung der Arbeitslosenquote und jährliches Wachstum des Bruttoinlandsprodukts zwischen 1992 und 2012 für Deutschland (*links*) beziehungsweise für die 6 Gründungsstaaten der EU (*rechts*) sowie jeweilige Regressionsgrade

### 2.1.1 Lineare Regression

Wir beginnen mit der einfachen linearen Regression:

**Definition 2.1** Im Modell der **einfachen linearen Regression** werden

$$Y_i = ax_i + b + \varepsilon_i, \quad i \in \{1, \dots, n\},$$

für gegebene Kovariablen  $x_1, \dots, x_n \in \mathbb{R}$  beobachtet. Hierbei sind die Beobachtungsfehler  $\varepsilon_1, \dots, \varepsilon_n$  zentrierte und unkorrelierte Zufallsvariablen ( $\mathbb{E}[\varepsilon_i] = 0$ ) mit endlicher Varianz  $\text{Var}(\varepsilon_i) = \sigma^2 > 0$ . Die Parameter  $a, b \in \mathbb{R}$  sind unbekannt und bestimmen die **Regressionsgerade**  $y = ax + b$ .

Aufgrund der Beobachtungsfehler  $\varepsilon_i$  ist der Zusammenhang zwischen den Kovariablen  $x_i$  und den Regressanden  $Y_i$  nicht deterministisch. Stattdessen beobachten wir eine zufällige Punktwolke um die Regressiongerade, die den linearen Zusammenhang möglichst gut beschreibt. Das Ziel ist die Schätzung der Parameter  $a$  und  $b$ . Der Parameter  $\sigma$  ist typischerweise nicht das Ziel der statistischen Inferenz und somit ein *Störparameter*.

**Beispiel 2.2 (Einfache lineare Regression, Okuns Gesetz)**  $Y_i$  bezeichnet das Wachstum des Bruttoinlandsprodukts von Deutschland im Jahr  $i$ . Die Kovariable  $x_i$  ist die Veränderung der Arbeitslosenquote im Vergleich zum Vorjahr. Unter Verwendung der Daten von 1992 bis 2012 aus den *World Development Indicators* der Weltbank erhalten wir als Regressionsgrade  $y = -1,080x + 1,338$ . Betrachten wir alle sechs Gründungsmitglieder der EU im gleichen Zeitraum, ergibt sich ganz ähnlich  $y = -1,075x + 1,819$ , siehe Abbildung 2.1. Der lineare Zusammenhang beider Größen ist als *Okuns Gesetz* bekannt.

Um die den Daten zugrunde liegenden Parameter  $a, b$  und damit die Regressionsgerade zu schätzen, werden wir die uns schon bekannte Maximum-Likelihood-Methode anwenden. Dafür nehmen wir an, dass  $\varepsilon_1, \dots, \varepsilon_n$  unabhängig und  $N(0, \sigma^2)$ -verteilt sind. Weil  $ax_i + b$  für alle  $i$  deterministisch in  $\mathbb{R}$  ist, gilt für die Beobachtungen in der einfachen linearen Regression

$$Y_i = ax_i + b + \varepsilon_i \sim N(ax_i + b, \sigma^2).$$

Das statistische Modell ist somit durch

$$\left( \mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \left( \bigotimes_{i=1}^n N(ax_i + b, \sigma^2) \right)_{a,b \in \mathbb{R}, \sigma > 0} \right)$$

gegeben. Es ergibt sich für  $y = (y_1, \dots, y_n) \in \mathbb{R}^n$  die Likelihood-Funktion

$$\begin{aligned} L(a, b, \sigma; y) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y_i - ax_i - b)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - ax_i - b)^2\right). \end{aligned}$$

Die Terme  $y_i - ax_i - b$  nennt man auch *Residuen*. Das Maximieren der Likelihood über  $a, b$  ist also äquivalent zum Minimieren der Summe der quadrierten Residuen (englisch: *residual sum of squares*, kurz: RSS). Auch wenn die Fehler nicht normalverteilt sind, kann diese Methode gute Ergebnisse erzielen.

**Methode 2.3 (Kleinste Quadrate)** In der einfachen linearen Regression sind die Kleinste-Quadrate-Schätzer (englisch: *least squares estimator*, kurz: LSE)  $\hat{a}, \hat{b}$  durch Minimierung der Summe der quadrierten Residuen gegeben:

$$(\hat{a}, \hat{b}) := \arg \min_{(a,b) \in \mathbb{R}^2} \sum_{i=1}^n (Y_i - ax_i - b)^2$$

Da die Kleinste-Quadrate-Schätzer nicht von  $\sigma^2$  abhängen, sind sie auch im Fall von unbekannter Fehlervarianz anwendbar. Zudem kann die Lösung dieses Minimierungsproblems explizit berechnet werden.

**Lemma 2.4** In der einfachen linearen Regression mit unabhängigen und  $N(0, \sigma^2)$ -verteilten Fehlern ist der Maximum-Likelihood-Schätzer gleich dem Kleinste-Quadrate-Schätzer. Zudem gilt

$$\hat{a} = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(Y_i - \bar{Y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \quad \text{und} \quad \hat{b} = \bar{Y}_n - \hat{a} \bar{x}_n$$

mit  $\bar{Y}_n := \frac{1}{n} \sum_{i=1}^n Y_i$  und  $\bar{x}_n := \frac{1}{n} \sum_{i=1}^n x_i$ , falls es  $i, j \in \{1, \dots, n\}$  gibt mit  $x_i \neq x_j$ .

**Beweis** Es bleibt festzustellen, dass wir durch Differenzieren in  $a$  und  $b$  folgende Normalgleichungen erhalten:

$$0 = \sum_{i=1}^n x_i(Y_i - ax_i - b) \quad \text{und} \quad 0 = \sum_{i=1}^n (Y_i - ax_i - b)$$

Man prüft leicht nach, dass diese Gleichungen durch  $\hat{a}$  und  $\hat{b}$  gelöst werden, sofern die Stichprobenvarianz der  $(x_i)$  nicht null ist, das heißt falls  $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 > 0$  gilt. Dies ist genau dann der Fall, wenn es  $x_i, x_j$  mit  $x_i \neq x_j$  gibt. Offensichtlich liegt bei  $(\hat{a}, \hat{b})$  ein Minimum des streng konvexen Kleinste-Quadrate-Kriteriums vor.  $\square$

Der Maximum-Likelihood-Ansatz liefert auch unter anderen Verteilungsannahmen an die  $(\varepsilon_i)$  sinnvolle Schätzer, siehe Aufgabe 2.1. Die Kleinste-Quadrate-Methode ist jedoch aufgrund ihrer guten Eigenschaften mit Abstand die populärste.

**Kurzbiografie (Carl Friedrich Gauß)** Carl Friedrich Gauß wurde 1777 in Braunschweig geboren. Früh galt er als Wunderknabe, was ihm die Gönnerschaft des Herzogs von Braunschweig einbrachte. Er unterstützte ihn vor allem finanziell, sodass er ab 1795 in Göttingen studieren konnte. Gauß beschäftigte sich mit Philologie, Mathematik, Physik und insbesondere mit Astronomie. An der Universität Helmstedt reichte er 1799 seine Doktorarbeit ein. Danach arbeitete er intensiv an seinem einflussreichen Lehrbuch zu höherer Arithmetik, der *Disquisitiones Arithmeticae*. 1801 gelang es ihm, den Zwergplaneten Ceres mithilfe seiner *Methode der kleinsten Quadrate* wiederaufzufinden. Wenig später wurde er Universitätsprofessor und Direktor der Sternwarte in Göttingen. Zahlreiche Methoden und Ideen sind nach ihm benannt, unter anderem das gaußsche Eliminationsverfahren zur Diagonalisierung und Invertierung von Matrizen, die Gauß-Verteilung (trotz reichlicher Vorarbeit von de Moivre, Laplace und Poisson) und die Methode der kleinsten Quadrate. 1855 starb Gauß in Göttingen.

Auch wenn kein linearer Zusammenhang vorliegt, kann die lineare Regression zur Untersuchung aller möglichen Zusammenhänge herangezogen werden: Der Abstand zur nächsten Autobahn (Regressor) als Einfluss auf die Anzahl der nächtlich geschlafenen Stunden (Regressand), der Betrag des Geldes auf dem Konto (Regressor) als Einfluss auf das Wohlbefinden (Regressand), die Anzahl der Sonnenstunden (Regressor) als Einfluss auf den Vitamin-D-Gehalt im Körper (Regressand) etc. Dabei muss man jedoch beachten, dass die Kalibrierung eines linearen Modells anhand von Daten, die keinem linearen Zusammenhang folgen, zu einem möglicherweise großen Modellfehler führen und gegebenenfalls falsche Schlussfolgerungen suggerieren, siehe Beispiel 2.9 sowie Aufgabe 2.3.

Bei der einfachen linearen Regression wird der Regressand durch eine Kovariable erklärt. Allgemeiner kann man den Einfluss mehrerer Kovariablen untersuchen.

**Definition 2.5** Bei  $k \geq 2$  Kovariablen  $x_j = (x_{1,j}, \dots, x_{n,j})^\top \in \mathbb{R}^n$ ,  $j = 1, \dots, k$ , und  $n$  Beobachtungen  $Y_i$  erhalten wir das **multiple lineare Regressionsmodell**

$$Y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{i,j} + \varepsilon_i, \quad i = 1, \dots, n,$$

wobei die Fehlerterme  $(\varepsilon_i)_{i=1, \dots, n}$  zentriert und unkorreliert sind mit  $0 < \text{Var}(\varepsilon_i) =: \sigma^2 < \infty$ . In Vektorschreibweise erhalten wir die lineare Gleichung

$$Y = X\beta + \varepsilon$$

mit

$$\begin{aligned} \textbf{Responsevektor} \quad Y &= (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n, \\ \textbf{Designmatrix} \quad X &:= \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} \end{pmatrix} \in \mathbb{R}^{n \times (k+1)}, \\ \textbf{Fehlervektor} \quad \varepsilon &:= (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n, \\ \textbf{Parametervektor} \quad \beta &:= (\beta_0, \dots, \beta_k)^\top \in \mathbb{R}^{k+1}. \end{aligned}$$

*Bemerkung 2.6* Wechselwirkungen zwischen zwei Kovariablen  $x_i$  und  $x_j$  werden in der Praxis oft durch Interaktionsterme  $x_{i,k} \cdot x_{i,j}$  in der Designmatrix modelliert.

Kategorielle Kovariablen können durch eine Menge von sogenannten *Dummy-Variablen*, das heisst  $\{0, 1\}$ -wertigen Variablen, kodiert werden, um nicht implizit eine (inadäquate) Metrisierung auf dem diskreten Wertebereich solcher Kovariablen zu erzeugen. Eine kategorielle Kovariable mit  $\ell$  möglichen Ausprägungen wird dabei durch  $(\ell - 1)$  viele  $\{0, 1\}$ -wertige Variablen repräsentiert. Die  $j$ -te Dummy-Variable kodiert das Ereignis, dass die Kategorie  $(j + 1)$  bei der zugehörigen Kovariablen vorliegt. Sind alle  $(\ell - 1)$  Indikatoren gleich null, so entspricht dies der (Referenz-) Kategorie 1 der zugehörigen kategoriellen Kovariablen. Wir werden dieses Vorgehen insbesondere bei der Varianzanalyse in Kapitel 2.3 verwenden, wo die *Faktoren* kategorielle Kovariablen sind.

Die Vektorschreibweise führt uns auf die allgemeine Form des linearen Modells. Dabei muss insbesondere die erste Spalte der Designmatrix nicht nur aus Einsen bestehen. Zudem werden Korrelationen zwischen den Fehlertermen  $\varepsilon_i$  zugelassen. Zunächst betrachten wir den klassischen Fall, in dem die Parameterdimension  $p$  kleiner als die Stichprobengröße  $n$  ist. Für den Fall  $p \geq n$  werden wir in Kapitel ?? eine Schätzmethode kennenlernen.

**Definition 2.7** Ein **lineares Modell** mit  $n$  reellwertigen Beobachtungen  $Y = (Y_1, \dots, Y_n)^\top$  und  $p$ -dimensionalem Parameter  $\beta \in \mathbb{R}^p$ ,  $p \leq n$ , besteht aus einer reellen Matrix  $X \in \mathbb{R}^{n \times p}$  von vollem Rang  $p$ , der **Designmatrix**, und einem Zufallsvektor  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$ , den **Fehler- oder Störgrößen**, mit  $\mathbb{E}[\varepsilon_i] = 0$  und positiv definiter Kovarianzmatrix  $\Sigma := \text{Cov}(\varepsilon) = (\text{Cov}(\varepsilon_i, \varepsilon_j))_{i,j=1, \dots, n}$ . Beobachtet wird eine Realisierung von

$$Y = X\beta + \varepsilon.$$

Wir sprechen vom **gewöhnlichen linearen Modell**, falls  $\Sigma = \sigma^2 E_n$  für ein Fehler-niveau  $\sigma > 0$  gilt.

*Bemerkung 2.8 (symmetrisch, positiv-definit)* Wir schreiben  $\Sigma > 0$ , falls  $\Sigma$  eine symmetrische, positiv-definite Matrix ist. Dann ist  $\Sigma = TDT^\top$  diagonalisierbar mit einer Diagonalmatrix  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ , Eigenwerten  $\lambda_i > 0$ ,  $i = 1, \dots, n$ , und einer Orthogonalmatrix  $T$ . Wir setzen  $\Sigma^{-1/2} := TD^{-1/2}T^\top$  mit  $D^{-1/2} := \text{diag}(\lambda_1^{-1/2}, \dots, \lambda_n^{-1/2})$  und erhalten

$$(\Sigma^{-1/2})^2 = \Sigma^{-1} \quad \text{und} \quad |\Sigma^{-1/2}v|^2 = \langle \Sigma^{-1}v, v \rangle.$$

Zusätzlich zur einfachen und multiplen Regression umfasst das lineare Modell weitere wichtige Beispiele.

*Beispiel 2.9 (Polynomiale Regression)* Wir beobachten für ein  $p \in \mathbb{N}$

$$Y_i = a_0 + a_1 x_i + a_2 x_i^2 + \dots + a_{p-1} x_i^{p-1} + \varepsilon_i, \quad i = 1, \dots, n.$$

Die Regressionsfunktion ist damit keine Gerade mehr, sondern ein Polynom vom Grad  $p-1$ . Die Koeffizienten des Polynoms bilden den unbekannten Parametervektor  $\beta = (a_0, \dots, a_{p-1})^\top$ . Es ergibt sich eine Designmatrix vom Vandermonde-Typ

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{p-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{p-1} \end{pmatrix}.$$

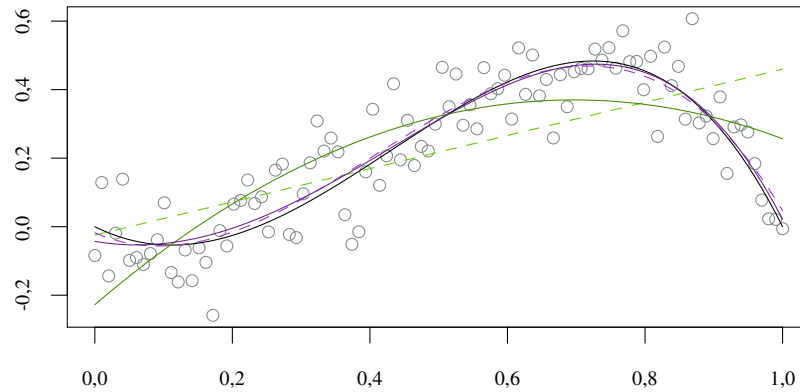
Die Matrix hat vollen Rang, sofern  $p$  der Designpunkte  $(x_i)$  verschieden sind, was über die sogenannte Vandermonde-Determinante leicht nachzuweisen ist. Abbildung 2.2 zeigt  $n = 100$  Beobachtungen, die durch ein Regressionspolynom vom Grad vier, äquidistante Designpunkte  $x_i = (i-1)/(n-1)$  und i.i.d. Beobachtungsfehler  $\varepsilon_i \sim N(0; 0, 1)$  erzeugt wurden.

*Beispiel 2.10 (Orthogonales Design)* Beobachten wir  $(x_i, Y_i)_{i=1, \dots, n}$  für reellwertige Kovariablen  $x_i \in \mathbb{R}$  und Regressanden  $Y_i$ , können wir das Regressionsmodell in der Form

$$Y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

mit einer unbekannten Regressionsfunktion  $f$  schreiben. In der polynomiellen Regression postulieren wir, dass  $f$  ein Polynom vom Grad  $p-1$  ist, und beschreiben  $f$  als Linearkombination der ersten  $p$  Monome  $(x^{k-1})_{k=1, \dots, p}$ . Analog können andere Basisfunktionen  $(\varphi_k)_{k=1, \dots, p}$  verwendet werden, um  $f(x) = \sum_{k=1}^p \beta_k \varphi_k(x)$  zu modellieren. Im Hinblick auf das Beobachtungsschema ist es nützlich, wenn die  $(\varphi_i)_{i=1, \dots, p}$  ein Orthonormalsystem bezüglich des *empirischen Skalarproduktes*

$$\langle f, g \rangle_n := \frac{1}{n} \sum_{i=1}^n f(x_i)g(x_i)$$



**Abb. 2.2** Beobachtungen  $(X_i, Y_i)_{i=1, \dots, n}$  aus einem polynomialen Regressionsmodell vom Grad 4 sowie die wahre Regressionsfunktion (schwarz) und geschätzte Regressionspolynome vom Grad eins (grün, gestrichelt), zwei (grün, durchgezogen), drei (violett, gestrichelt), vier (violett, durchgezogen)

bilden. In diesem Fall besitzt die Designmatrix  $X = (\varphi_j(x_i))_{i=1, \dots, n, j=1, \dots, p}$  die Eigenschaft  $X^T X = nE_p$ , weshalb wir von *orthogonalem Design* sprechen. Beispielsweise können die Monome mittels Gram-Schmidt-Verfahren bezüglich  $\langle \cdot, \cdot \rangle_n$  orthogonalisiert werden. Es sei bemerkt, dass für äquidistantes Design das empirische Skalarprodukt für  $n \rightarrow \infty$  gegen das  $L^2$ -Skalarprodukt konvergiert, siehe hierzu Beispiel 3.13.

### 2.1.2 Schätzen im linearen Modell

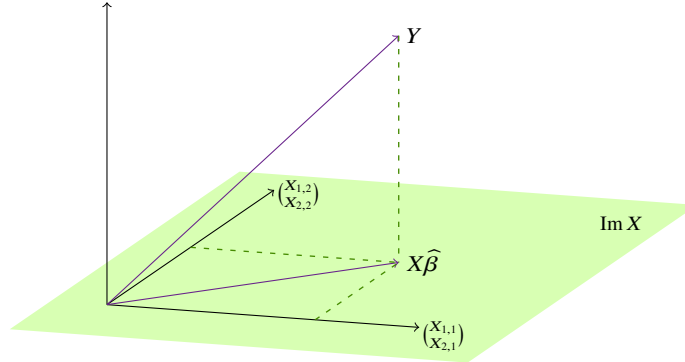
Übertragen wir die Kleinste-Quadrate-Methode (siehe Methode 2.3) von der einfachen linearen Regression auf den allgemeinen Fall, erhalten wir folgendes Verfahren zur Schätzung des Parametervektors  $\beta$  im linearen Modell:

**Methode 2.11 (Kleinste-Quadrate-Schätzer)** Der **gewichtete Kleinste-Quadrate-Schätzer**  $\hat{\beta}$  von  $\beta$  minimiert den gewichteten euklidischen Abstand zwischen Beobachtungen und Modellvorhersage:

$$|\Sigma^{-1/2}(Y - X\hat{\beta})|^2 = \inf_{b \in \mathbb{R}^p} |\Sigma^{-1/2}(Y - Xb)|^2$$

Im gewöhnlichen Fall  $\Sigma = \sigma^2 E_n$  ergibt sich der **gewöhnliche Kleinste-Quadrate-Schätzer** (englisch: *ordinary least squares*, kurz: OLS)  $\hat{\beta}$  mit

$$|Y - X\hat{\beta}|^2 = \inf_{b \in \mathbb{R}^p} |Y - Xb|^2,$$



**Abb. 2.3** Geometrische Interpretation des Kleinst-Quadrate-Schätzers mit  $n = 3$  und  $p = 2$

der unabhängig von der Kenntnis von  $\sigma^2$  ist.

*Bemerkung 2.12 (Gewichteter Kleinst-Quadrate-Schätzer)* In der einfachen linearen Regression hatten wir den Kleinst-Quadrate-Schätzer mit normalverteilten, unabhängigen und identisch verteilten Fehlern hergeleitet. Das allgemeine lineare Modell können wir hierauf zurückführen, indem wir die beobachteten Daten entsprechend der Kovarianzmatrix  $\Sigma$  gewichten, genauer betrachten wir  $\Sigma^{-1/2}Y$ . Für die entsprechend gewichteten Fehler  $\Sigma^{-1/2}\varepsilon$  gilt nämlich

$$\text{Cov}(\Sigma^{-1/2}\varepsilon) = \Sigma^{-1/2} \text{Cov}(\varepsilon)(\Sigma^{-1/2})^\top = \Sigma^{-1/2}\Sigma\Sigma^{-1/2} = E_n.$$

Aus  $\Sigma^{-1/2}\varepsilon = \Sigma^{-1/2}(Y - X\beta)$  folgt der Ansatz des gewichteten Kleinst-Quadrate-Schätzers.

*Beispiel 2.13 (Polynomiale Regression)* Im Modell aus Beispiel 2.9 und Abbildung 2.2 wenden wir nun die Kleinst-Quadrate-Methode zur Schätzung der unbekannten Koeffizienten des zugrunde liegenden Regressionspolynoms an. Anhand der Beobachtungen aus Abbildung 2.2 kalibrieren wir polynomiale Regressionsmodelle der Grade 1, 2, 3 und 4. Da in den ersten drei Fällen das zur Schätzung verwendete Modell nicht mit dem wahren Modell, das die Daten erzeugt hat, übereinstimmt, weisen diese Schätzungen einen Modellfehler auf. Dieser zeigt sich insbesondere für die Grade eins und zwei in einer deutlichen Abweichung der geschätzten von der wahren Regressionsfunktion. Andererseits scheint bereits ein Polynom vom Grad drei die Daten gut zu beschreiben, da der zusätzliche vierte Grad kaum zu Änderungen führt.

Wir betrachten das gewöhnliche lineare Modell mit  $\Sigma = E_n$  und einem zweidimensionalen Parameter  $\beta \in \mathbb{R}^2$ , das heißt  $p = 2$ . Die Designmatrix  $X$  ist dann eine  $(n \times 2)$ -Matrix und wir beobachten einen Punkt  $Y \in \mathbb{R}^n$ . Der Kleinst-Quadrate-Schätzer gibt nun den Wert  $b = \hat{\beta}$  an, an dem der euklidische Abstand zwischen  $Y$  und der Ebene  $\text{Im } X = \{Xb : b \in \mathbb{R}^2\}$  minimal ist. Folglich ist  $X\hat{\beta}$  die *Orthogonalprojektion* von  $Y$  auf  $\text{Im } X$ , siehe Abbildung 2.3.

Diese geometrische Anschauung lässt sich auch ganz allgemein formulieren und führt zu einer expliziten Lösung des Minimierungsproblems.

**Lemma 2.14** *Es seien  $\Sigma > 0$  und  $X$  von vollem Rang. Setze*

$$X_\Sigma := \Sigma^{-1/2}X \quad \text{und} \quad \Pi_{X_\Sigma} := X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1}X_\Sigma^\top.$$

*Dann ist  $\Pi_{X_\Sigma}$  die Orthogonalprojektion von  $\mathbb{R}^n$  auf den Bildraum  $\text{Im}(X_\Sigma) = \{X_\Sigma b : b \in \mathbb{R}^p\}$  und für den Kleinst-Quadrate-Schätzer gilt*

$$\widehat{\beta} = X_\Sigma^{-1}(\Pi_{X_\Sigma} \Sigma^{-1/2}Y) = (X^\top \Sigma^{-1}X)^{-1}X^\top \Sigma^{-1}Y.$$

*Dabei ist  $X_\Sigma^{-1}(v)$  für  $v \in \text{Im}(X_\Sigma)$  das eindeutig bestimmte Urbild von  $v$  unter  $X_\Sigma$  (für  $p < n$  ist  $X_\Sigma$  keine invertierbare Matrix). Insbesondere existiert der Kleinst-Quadrate-Schätzer, ist eindeutig und erwartungstreu, und es gilt  $X_\Sigma \widehat{\beta} = \Pi_{X_\Sigma}(\Sigma^{-1/2}Y)$ .*

**Beweis** Die Positivität und Symmetrie von  $\Sigma$  liefert

$$\Sigma^{-1/2} = (\Sigma^{-1/2})^\top \quad \text{und} \quad X_\Sigma^\top = (\Sigma^{-1/2}X)^\top = X^\top (\Sigma^{-1/2})^\top = X^\top \Sigma^{-1/2}. \quad (2.1)$$

Wir zeigen zuerst, dass  $X_\Sigma^\top X_\Sigma = X^\top \Sigma^{-1}X$  invertierbar und somit  $\Pi_{X_\Sigma}$  wohldefiniert ist: Für jedes  $v \in \mathbb{R}^p$  mit  $X^\top \Sigma^{-1}Xv = 0$  folgt aus (2.1)

$$0 = v^\top X^\top \Sigma^{-1}Xv = |\Sigma^{-1/2}Xv|^2.$$

Da  $\Sigma^{-1/2}$  vollen Rang hat, muss dann  $|Xv| = 0$  gelten. Aus dem vollen Rang von  $X$  folgt wiederum  $v = 0$ . Also besteht der Kern von  $X_\Sigma^\top X_\Sigma$  nur aus dem Nullvektor, und  $X_\Sigma^\top X_\Sigma$  ist invertierbar.

Wir setzen nun  $\Pi_{X_\Sigma} := X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1}X_\Sigma^\top$  und  $w = \Pi_{X_\Sigma}v$  für ein  $v \in \mathbb{R}^n$ . Dann folgt  $w \in \text{Im}(X_\Sigma)$  und im Fall  $v = X_\Sigma u$  durch Einsetzen  $w = \Pi_{X_\Sigma}X_\Sigma u = v$ , sodass  $\Pi_{X_\Sigma}$  eine Projektion auf  $\text{Im}(X_\Sigma)$  ist. Zudem ist  $\Pi_{X_\Sigma}$  selbstadjungiert (symmetrisch) wegen

$$\begin{aligned} ((X_\Sigma^\top X_\Sigma)^{-1})^\top &= ((X^\top \Sigma^{-1}X)^{-1})^\top = ((X^\top \Sigma^{-1}X)^\top)^{-1} \\ &= (X^\top (\Sigma^{-1})^\top X)^{-1} = (X^\top \Sigma^{-1}X)^{-1} = (X_\Sigma^\top X_\Sigma)^{-1}. \end{aligned} \quad (2.2)$$

Somit ist  $\Pi_{X_\Sigma}$  sogar eine Orthogonalprojektion:

$$\forall u \in \mathbb{R}^n, \forall w \in \text{Im}(X_\Sigma) : \langle u - \Pi_{X_\Sigma}u, w \rangle = \langle u, w \rangle - \langle u, \Pi_{X_\Sigma}w \rangle = 0$$

Aus der Eigenschaft  $\widehat{\beta} = \arg \min_b |\Sigma^{-1/2}(Y - Xb)|^2$  folgt, dass  $\widehat{\beta}$  die beste Approximation von  $\Sigma^{-1/2}Y$  durch  $X_\Sigma b$  liefert. Diese ist durch die Orthogonalprojektionseigenschaft

$$\Pi_{X_\Sigma} \Sigma^{-1/2}Y = X_\Sigma \widehat{\beta}$$

bestimmt. Es folgt

$$X_{\Sigma}^{\top} \Pi_{X_{\Sigma}} \Sigma^{-1/2} Y = (X_{\Sigma}^{\top} X_{\Sigma}) \widehat{\beta} \Rightarrow (X_{\Sigma}^{\top} X_{\Sigma})^{-1} X_{\Sigma}^{\top} \Sigma^{-1} Y = \widehat{\beta}.$$

Schließlich folgt aus der Linearität des Erwartungswerts und  $\mathbb{E}[\varepsilon] = 0$

$$\mathbb{E}[\widehat{\beta}] = \mathbb{E}[(X_{\Sigma}^{\top} X_{\Sigma})^{-1} X_{\Sigma}^{\top} \Sigma^{-1} (X\beta + \varepsilon)] = \beta + 0 = \beta.$$

Damit ist  $\widehat{\beta}$  erwartungstreu. □

*Bemerkung 2.15*

- Im gewöhnlichen linearen Modell gilt  $\widehat{\beta} = (X^{\top} X)^{-1} X^{\top} Y$ , sodass der Kleinste-Quadrate-Schätzer unabhängig vom unbekannten Parameter  $\sigma > 0$  ist.
- $X_{\Sigma}^{\dagger} := (X_{\Sigma}^{\top} X_{\Sigma})^{-1} X_{\Sigma}^{\top}$  heißt auch *Moore-Penrose-(Pseudo-)Inverse* von  $X_{\Sigma}$ . Die Bezeichnung Pseudoinverse ist motiviert durch die Eigenschaften  $X_{\Sigma}^{\dagger} X_{\Sigma} = E_p$  und  $X_{\Sigma} X_{\Sigma}^{\dagger}|_{\text{Im}(X_{\Sigma})} = E_n|_{\text{Im}(X_{\Sigma})}$ . Insbesondere erhalten wir  $\widehat{\beta} = X_{\Sigma}^{\dagger} \Sigma^{-1/2} Y$  bzw. die Vereinfachung  $\widehat{\beta} = X^{\dagger} Y$  im gewöhnlichen linearen Modell.

Der folgende zentrale Satz der Regressionsanalyse zeigt, dass der Kleinste-Quadrate-Schätzer optimal ist, wenn wir uns auf Schätzer beschränken, die linear in den Daten  $Y$  und erwartungstreu sind.

**Satz 2.16 (Gauß-Markov)** *Im linearen Modell gilt:*

- Der Parameter  $\rho = \langle \beta, v \rangle$  für ein  $v \in \mathbb{R}^p$  wird von  $\widehat{\rho} = \langle \widehat{\beta}, v \rangle$  erwartungstreu geschätzt, und  $\widehat{\rho}$  besitzt unter allen linearen erwartungstreuen Schätzern die minimale Varianz  $\text{Var}(\widehat{\rho}) = |X_{\Sigma} (X_{\Sigma}^{\top} X_{\Sigma})^{-1} v|^2$ .
- Der Kleinste-Quadrate-Schätzer  $\widehat{\beta}$  besitzt unter allen linearen erwartungstreuen Schätzern von  $\beta$  minimale Kovarianzmatrix, nämlich  $\text{Cov}(\widehat{\beta}) = (X_{\Sigma}^{\top} X_{\Sigma})^{-1}$ , das heißt für alle linearen, erwartungstreuen Schätzer  $\widetilde{\beta}$  ist  $\text{Cov}(\widetilde{\beta}) - \text{Cov}(\widehat{\beta}) \geq 0$  eine positiv semi-definite Matrix.

**Beweis** (i) Die Linearität ist klar und aus dem vorangegangenen Lemma folgt, dass  $\widehat{\rho}$  erwartungstreu ist. Für die Varianz ergibt sich

$$\begin{aligned} \text{Var}(\widehat{\rho}) &= \mathbb{E}[\langle \widehat{\beta} - \beta, v \rangle^2] \\ &= \mathbb{E}[\langle (X^{\top} \Sigma^{-1} X)^{-1} X^{\top} \Sigma^{-1} \varepsilon, v \rangle^2] \\ &= \mathbb{E}[\langle \varepsilon, \Sigma^{-1} X (X^{\top} \Sigma^{-1} X)^{-1} v \rangle^2] \\ &= v^{\top} (X^{\top} \Sigma^{-1} X)^{-1} X^{\top} \Sigma^{-1} \Sigma \Sigma^{-1} X (X^{\top} \Sigma^{-1} X)^{-1} v \\ &= |X_{\Sigma} (X_{\Sigma}^{\top} X_{\Sigma})^{-1} v|^2. \end{aligned}$$

Sei nun  $\widetilde{\rho}$  ein beliebiger linearer Schätzer von  $\rho$ . Dann gibt es ein  $w \in \mathbb{R}^n$ , sodass  $\widetilde{\rho} = \langle Y, w \rangle$  (Satz von Riesz). Dies impliziert, dass für alle  $\beta \in \mathbb{R}^p$

$$\mathbb{E}[\langle Y, w \rangle] = \rho \Rightarrow \mathbb{E}[\langle Y, w \rangle] = \langle X\beta, w \rangle = \langle \beta, v \rangle \Rightarrow \langle X^{\top} w - v, \beta \rangle = 0$$

und somit  $v = X^{\top} w = X_{\Sigma}^{\top} \Sigma^{1/2} w$ . Da  $\Pi_{X_{\Sigma}}$  eine Orthogonalprojektion ist, erhalten wir mit dem Satz von Pythagoras

$$\begin{aligned}\text{Var}(\tilde{\rho}) &= \mathbb{E}[\langle \varepsilon, w \rangle^2] = \mathbb{E}[w^\top \varepsilon \varepsilon^\top w] \\ &= w^\top \Sigma w = |\Sigma^{1/2} w|^2 = |\Pi_{X_\Sigma}(\Sigma^{1/2} w)|^2 + |(E_n - \Pi_{X_\Sigma})(\Sigma^{1/2} w)|^2.\end{aligned}$$

Damit gilt

$$\begin{aligned}\text{Var}(\tilde{\rho}) &\geq |\Pi_{X_\Sigma}(\Sigma^{1/2} w)|^2 = |X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top w|^2 \\ &= |X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1} v|^2 = \text{Var}(\hat{\rho}).\end{aligned}$$

(ii) Sei  $\tilde{\beta}$  ein linearer, erwartungstreuer Schätzer. Dann ist  $\mathbb{Cov}(\tilde{\beta}) - \mathbb{Cov}(\hat{\beta})$  genau dann eine positiv semi-definite Matrix, wenn

$$\forall v \in \mathbb{R}^p : v^\top (\mathbb{Cov}(\tilde{\beta}) - \mathbb{Cov}(\hat{\beta})) v \geq 0.$$

Sei  $v \in \mathbb{R}^p$ . Nach Annahme sind  $\langle v, \tilde{\beta} \rangle$  und  $\langle v, \hat{\beta} \rangle$  lineare, erwartungstreue Schätzer und aus (i) folgt daher

$$0 \leq \text{Var}(\langle v, \tilde{\beta} \rangle) - \text{Var}(\langle v, \hat{\beta} \rangle) = v^\top (\mathbb{Cov}(\tilde{\beta}) - \mathbb{Cov}(\hat{\beta})) v.$$

Weiter gilt wegen (i) für beliebiges  $v \in \mathbb{R}^p$

$$\begin{aligned}v^\top \mathbb{Cov}(\hat{\beta}) v &= \text{Var}(\langle \hat{\beta}, v \rangle) \\ &= |X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1} v|^2 \\ &= v^\top (X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top X_\Sigma (X_\Sigma^\top X_\Sigma)^{-1} v \\ &= v^\top (X_\Sigma^\top X_\Sigma)^{-1} v,\end{aligned}$$

woraus  $\mathbb{Cov}(\hat{\beta}) = (X_\Sigma^\top X_\Sigma)^{-1}$  folgt.  $\square$

**Kurzbiografie (Andrey Andreyevich Markov)** Andrey Andreyevich Markov wurde 1856 in Rjasan, rund 200 km südöstlich von Moskau, geboren und wuchs in St. Petersburg auf. Er studierte an der St. Petersburger Universität Mathematik und Physik, wobei er unter anderem Vorlesungen von Tschebyscheff besuchte. Kurz nach seiner Promotion wurde er Professor und in die russische Akademie der Wissenschaften gewählt. Markov leistete wichtige Beiträge zur Stochastik und Analysis und insbesondere zu stochastischen Prozessen, wobei etliche Resultate nach ihm benannt wurden, beispielsweise die Markov-Eigenschaft, die Markov-Ungleichung, der Markov-Prozess und der *Satz von Gauß-Markov*. Markov starb 1922.

Da wir das Funktional  $\rho := \rho(\beta) := \langle v, \beta \rangle$  mit  $\hat{\rho} = \rho(\hat{\beta})$  durch Einsetzen von  $\hat{\beta}$  schätzen, nennt man  $\hat{\rho}$  *plugin-Schätzer*. Eine typische Anwendung sind Vorhersagen. Aufgrund des Schätzers  $\hat{\beta}$  von  $\beta$  können wir für neue Kovariablen  $v := (x_{n+1,1}, \dots, x_{n+1,p})^\top$  die zugehörige Beobachtung  $Y_{n+1} = v^\top \beta + \varepsilon_{n+1}$  vorhersagen, indem wir  $\langle v, \hat{\beta} \rangle$  berechnen. Eine weitere wichtige Anwendung ist die Schätzung einzelner Koeffizienten des Vektors  $\beta$ .

*Bemerkung 2.17 (BLUE, GLS)*

1. Man sagt, dass der Schätzer  $\hat{\rho}$  im Satz von Gauß-Markov **besten linearer erwartungstreuer Schätzer** (englisch: *best linear unbiased estimator*, kurz: BLUE) ist. Eingeschränkt auf lineare, erwartungstreue Schätzer ist der Kleinste-Quadrate-Schätzer damit minimax bezüglich quadratischem Verlust, siehe Definition 1.25. Ob es einen besseren nichtlinearen Schätzer geben kann, werden wir in Kapitel ?? beantworten.
2. Der Kleinste-Quadrate-Schätzer mit allgemeiner Kovarianzmatrix  $\Sigma$  wird im Englischen zur Abhebung vom gewöhnlichen Kleinste-Quadrate-Schätzer (OLS) auch verallgemeinerter Kleinste-Quadrate-Schätzer (englisch: *generalized least squares*, kurz: GLS) genannt und wurde erstmals von Alexander Aitken im Jahre 1934 beschrieben. Aitken hatte den Satz von Gauß-Markov, der sich eigentlich nur mit dem gewöhnlichen Fall beschäftigte, auf das allgemeine Modell übertragen.
3. Wegen  $\mathbb{E}[|Z|^2] = \sum_{i=1, \dots, p} \mathbb{E}[Z_i^2] = \text{tr}(\mathbb{E}[ZZ^\top])$  für Zufallsvektoren  $Z \in \mathbb{R}^p$  und mit der Spur  $\text{tr}(\cdot)$  einer Matrix ist der mittlere quadratische Fehler von  $\hat{\beta}$  gegeben durch  $\mathbb{E}[|\hat{\beta} - \beta|^2] = \text{tr}((X_\Sigma^\top X_\Sigma)^{-1})$ . Für ein gewöhnliches lineares Modell mit orthogonalem Design (Beispiel 2.10) gilt  $X_\Sigma^\top X_\Sigma = X^\top \Sigma^{-1} X = \frac{n}{\sigma^2} E_p$  und daher

$$\mathbb{E}[|\hat{\beta} - \beta|^2] = \frac{\sigma^2 p}{n}. \quad (2.3)$$

Die Abhängigkeit des Fehlers vom Stichprobenumfang  $n$ , der Parameterdimension  $p$  und dem Rauschniveau  $\sigma$  ist hier besonders offenkundig.

Im Spezialfall des gewöhnlichen linearen Modells ist es von großem Interesse, das *Rauschniveau*  $\sigma^2$  zu schätzen. Dieses ist der einzige unbekannte Wert in der Fehlerformel (2.3) und wird uns die Konstruktion von Tests und Konfidenzbereichen ermöglichen.

**Lemma 2.18** *Im gewöhnlichen linearen Modell mit  $\sigma > 0$  und Kleinste-Quadrate-Schätzer  $\hat{\beta}$  gilt  $X\hat{\beta} = \Pi_X Y$  (mit  $\Pi_X := \Pi_{X_{E_n}}$ ).  $R := Y - X\hat{\beta}$  bezeichne den Vektor der **Residuen**. Dann ist für  $p < n$  die Stichprobenvarianz*

$$\hat{\sigma}^2 := \frac{|R|^2}{n-p} = \frac{|(E_n - \Pi_X)Y|^2}{n-p}$$

*ein erwartungstreuer Schätzer von  $\sigma^2$ .*

**Beweis**  $X\hat{\beta} = \Pi_X Y$  folgt aus Lemma 2.14. Einsetzen zeigt  $\mathbb{E}[|Y - X\hat{\beta}|^2] = \mathbb{E}[|Y - \Pi_X Y|^2] = \mathbb{E}[|(E_n - \Pi_X)\varepsilon|^2]$ . Mithilfe der Spur und Eigenschaften der Orthogonalprojektion  $E_n - \Pi_X$  vom Rang  $n - p$  berechnen wir

$$\mathbb{E}[|(E_n - \Pi_X)\varepsilon|^2] = \text{tr}((E_n - \Pi_X)\mathbb{E}[\varepsilon\varepsilon^\top](E_n - \Pi_X)) = \sigma^2 \text{tr}(E_n - \Pi_X) = \sigma^2(n-p),$$

was die Behauptung impliziert.  $\square$

Eine alternative Normierung der Stichprobenvarianz ergibt sich aus dem Maximum-Likelihood-Ansatz im normalverteilten gewöhnlichen linearen Modell. Dieser

liefert  $\hat{\sigma}_{ML}^2 = |R|^2/n$  als Schätzer für  $\sigma^2$  (siehe Aufgabe 2.6). In der Praxis wird der erwartungstreue Schätzer  $\hat{\sigma}^2$  häufig bevorzugt, wobei dieser jedoch eine größere Varianz als  $\hat{\sigma}_{ML}^2$  aufweist.

Nachdem uns die Maximum-Likelihood-Methode auf den Kleinste-Quadrate-Schätzer geführt hat, soll nun der Bayes-Ansatz für das lineare Modell untersucht werden.

**Satz 2.19 (Bayes-Schätzer im linearen Modell)** *Im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$  und bekanntem  $\sigma > 0$  folge  $\beta \in \mathbb{R}^p$  der a-priori-Verteilung*

$$\beta \sim N(m, \sigma^2 M)$$

*mit  $m \in \mathbb{R}^p$  und symmetrischer, positiv definiter Matrix  $M \in \mathbb{R}^{p \times p}$ . Dann ist die a-posteriori-Verteilung von  $\beta$ , gegeben eine realisierte Beobachtung  $y \in \mathbb{R}^n$ , wiederum normalverteilt:*

$$\beta|Y=y \sim N(\mu_y, \sigma^2 \Sigma_y) \quad \text{mit} \quad \Sigma_y = (X^\top X + M^{-1})^{-1}, \quad \mu_y = \Sigma_y (X^\top y + M^{-1}m)$$

*Insbesondere ist der Bayes-Schätzer bezüglich quadratischem Verlust gegeben durch*

$$\hat{\beta}_{\text{Bayes}} = (X^\top X + M^{-1})^{-1} (X^\top Y + M^{-1}m).$$

**Beweis** Für die a-posteriori-Dichte an der Stelle  $t \in \mathbb{R}^p$  gilt

$$\begin{aligned} f^{\beta|Y=y}(t) &\propto f^{Y|\beta=t}(y) f_\beta(t) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}(y - Xt)^\top (y - Xt)\right) \exp\left(-\frac{1}{2\sigma^2}(t - m)^\top M^{-1}(t - m)\right) \\ &\propto \exp\left(\frac{1}{\sigma^2} t^\top X^\top y - \frac{1}{2\sigma^2} t^\top X^\top X t - \frac{1}{2\sigma^2} t^\top M^{-1} t + \frac{1}{\sigma^2} t^\top M^{-1} m\right) \\ &= \exp\left(\frac{1}{\sigma^2} t^\top \underbrace{(X^\top y + M^{-1}m)}_{=\Sigma_y^{-1}\mu_y} - \frac{1}{2\sigma^2} t^\top \underbrace{(X^\top X + M^{-1})}_{=\Sigma_y^{-1}} t\right) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} (t^\top \Sigma_y^{-1} t - 2t^\top \Sigma_y^{-1} \mu_y + \mu_y^\top \Sigma_y^{-1} \mu_y)\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} (t - \mu_y)^\top \Sigma_y^{-1} (t - \mu_y)\right). \end{aligned}$$

Daher ist  $\beta$ , gegeben  $Y = y$ , normalverteilt mit der Kovarianzmatrix  $\sigma^2 \Sigma_y$  und dem Erwartungswert  $\mu_y$ . Korollar 1.33 liefert den Rest der Behauptung.  $\square$

**Bemerkung 2.20 (Mehrstufiges Bayes-Modell)** Indem wir den Parameter  $\sigma^2$  mit einer a-priori-Verteilung versehen und  $\beta$  gemäß einer von  $\sigma$  abhängigen a-priori-Verteilung wählen, erhalten wir ein mehrstufiges Bayes-Modell. Da besonders konjugierte Verteilungsklassen von Interesse sind, wird hierzu oft die *inverse Gamma-Verteilung* verwendet: Ist  $Z \sim \Gamma(a, b)$ , so ist  $1/Z \sim \text{IG}(a, b)$  invers gammaverteilt

mit Parametern  $a, b > 0$  und Lebesgue-Dichte

$$f_{a,b}(x) = \frac{b^a}{\Gamma(a)} x^{-(a+1)} e^{-b/x} \mathbb{1}_{(0,\infty)}(x), \quad x \in \mathbb{R}.$$

Skalieren wir die Varianz des normalverteilten  $\beta$  mit  $\sigma^2$ , erhalten wir das Bayes-Modell

$$Y|\beta, \sigma^2 \sim N(X\beta, \sigma^2 E), \quad \beta|\sigma^2 \sim N(m, \sigma^2 M), \quad \sigma^2 \sim \text{IG}(a, b).$$

Die gemeinsame Verteilung von  $(\beta, \sigma^2) \sim \text{NIG}(m, M, a, b)$  wird *normal-inverse Gammaverteilung* genannt und besitzt die Dichte

$$f(\beta, \sigma^2) = \frac{C}{(\sigma^2)^{p/2+a+1}} \exp\left(-\frac{1}{2\sigma^2}((\beta - m)^\top M^{-1}(\beta - m) + 2b)\right) \quad \text{mit}$$

$$C = \frac{b^a}{(2\pi)^{p/2} |M|^{1/2} \Gamma(a)}, \quad \beta \in \mathbb{R}^p, \sigma^2 > 0.$$

In diesem Modell ist die a-posteriori-Verteilung von  $\sigma^2$ , gegeben  $\beta$  und  $Y$ , gegeben durch  $\sigma^2|\beta, Y \sim \text{IG}(a', b')$  mit  $a' = a + \frac{n}{2} + \frac{p}{2}$  und

$$b' = b + \frac{1}{2}(Y - X\beta)^\top (Y - X\beta) + \frac{1}{2}(\beta - m)^\top M^{-1}(\beta - m).$$

Die a-posteriori-Verteilung von  $(\beta, \sigma^2)$ , gegeben  $Y$ , ist  $(\beta, \sigma^2)|Y \sim \text{NIG}(\tilde{m}, \tilde{M}, \tilde{a}, \tilde{b})$  mit den Parametern

$$\tilde{M} = (X^\top X + M^{-1})^{-1}, \quad \tilde{m} = \tilde{M}(M^{-1}m + X^\top y),$$

$$\tilde{a} = a + \frac{n}{2}, \quad \tilde{b} = b + \frac{1}{2}(Y^\top Y + m^\top M^{-1}m - \tilde{m}^\top \tilde{M}^{-1}\tilde{m}),$$

siehe Fahrmeir et al. (2009, Kapitel 3.5).

In einem Spezialfall von Satz 2.19 erhalten wir eine weitere Darstellung des Bayes-Schätzers.

**Korollar 2.21** *Unter den Voraussetzungen von Satz 2.19 mit  $m = 0$  und  $M = \lambda^{-1}E_p$  für ein  $\lambda > 0$  gilt für den Bayes-Schätzer unter quadratischem Verlust*

$$\hat{\beta}_{\text{Bayes}} = \arg \min_{\beta \in \mathbb{R}^p} (|Y - X\beta|^2 + \lambda|\beta|^2).$$

**Beweis** Im Spezialfall  $m = 0$  und  $M = \lambda^{-1}E_p$  folgt aus Satz 2.19

$$\hat{\beta}_{\text{Bayes}} = (X^\top X + \lambda E_p)^{-1} X^\top Y.$$

Andererseits gilt

$$\arg \min_{\beta} (|Y - X\beta|^2 + \lambda \beta^\top \beta) = \arg \min_{\beta} (-2Y^\top X\beta + \beta^\top (X^\top X + \lambda E_p)\beta).$$

Null setzen des Differenzials in  $\beta$  liefert  $0 = -2Y^\top X + 2\beta^\top (X^\top X + \lambda E_p)$ , sodass aus der Positivität und Symmetrie von  $X^\top X + \lambda E_p$  die Behauptung folgt.  $\square$

Der Bayes-Ansatz führt uns also zu einer neuen Schätzmethode im linearen Modell:

**Methode 2.22 (Ridge-Regression)** Im linearen Modell  $Y = X\beta + \varepsilon$  ist der **Ridge-Regressionsschätzer** oder **Shrinkage-Schätzer** mit Penalisierung  $\lambda > 0$  definiert als

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta \in \mathbb{R}^p} (|Y - X\beta|^2 + \lambda |\beta|^2).$$

Der Strafterm  $\lambda |\beta|^2$  führt zu Lösungen des Minimierungsproblems, die kleine Parametervektoren  $\beta$  bevorzugen, und daher zur Bezeichnung *shrinkage* („Schrumpfung“) führt. Diesen Effekt sieht man auch direkt am Bayes-Ansatz in Satz 2.19, da dort eine a-priori-Verteilung verwendet wird, die um  $\beta = 0$  zentriert ist. Um die Bezeichnung *ridge* zu verstehen, erinnern wir uns daran, dass die ursprüngliche Motivation des Kleinste-Quadrate-Schätzers das Maximum-Likelihood-Prinzip war. Führen nun mehrere Parameterwahlen  $\beta$  zu vergleichbar großen Likelihoods, ähnelt die Kontur der Likelihood-Funktion einem Bergkamm (englisch: *ridge*). Das Finden des Maximums oder äquivalent des Minimums der quadrierten Residuen (in einem langen Tal) ist numerisch schwierig. Wenn wir durch Hinzufügen des strikt konvexen Strafterms das Tal an den Seiten anheben, entsteht ein leichter zu findendes globales Minimum.

Abbildung 2.4 illustriert diesen Sachverhalt, wobei wir für die Simulation folgendes Modell benutzt haben:  $X \in \mathbb{R}^{10 \times 2}$ ,  $X_{i,1} = 1$ ,  $X_{j,2} \sim U([0, 1])$  für alle  $i, j = 1, \dots, 10$ ,  $\beta = (2, 1/2)^\top$ ,  $\varepsilon_i \sim N(0, 1/10)$  mit  $\Sigma = 0$ ,  $1 \cdot E_{10}$  und  $\lambda = 15$ .

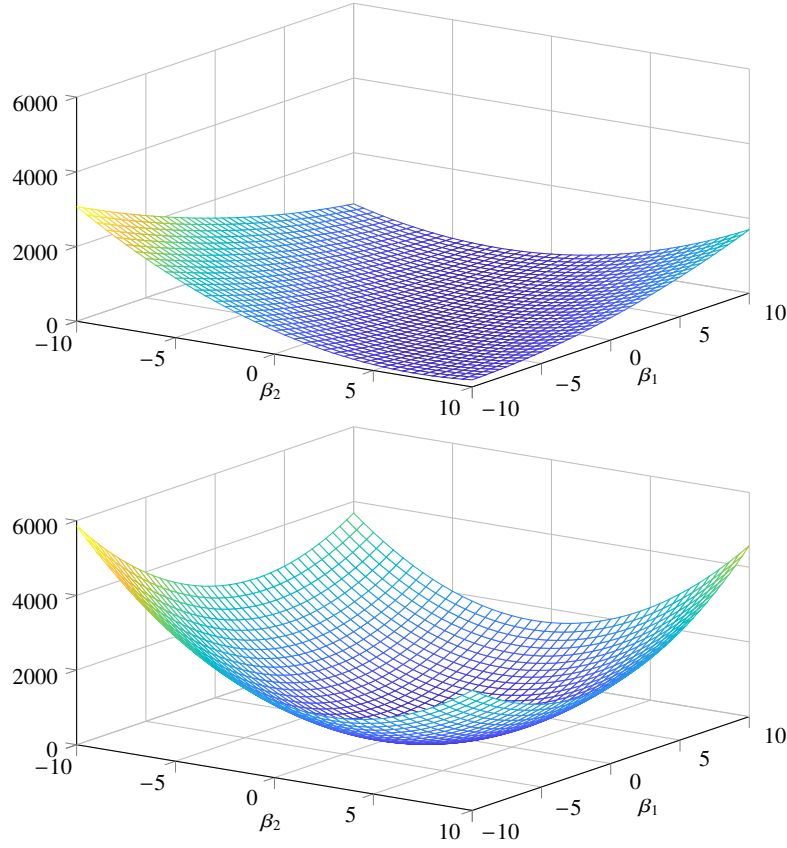
Den Einfluss des Strafterms wollen wir im Spezialfall von orthogonalem Design, das heißt für  $X^\top X = nE_p$  und  $p \leq n$ , genauer untersuchen, siehe auch Aufgabe 2.9 zum Vergleich zwischen Ridge-Regressionsschätzer und dem Kleinste-Quadrate-Schätzer.

**Lemma 2.23** Im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$ ,  $\sigma > 0$  und  $X^\top X = nE_p$  gilt für den Ridge-Regressionsschätzer mit Penalisierungsparameter  $\lambda > 0$

$$\mathbb{E}[|\hat{\beta}_{\text{ridge}} - \beta|^2] = \frac{|\beta|^2}{(1 + n/\lambda)^2} + \frac{1}{(1 + \lambda/n)^2} \frac{\sigma^2 p}{n}.$$

**Beweis** Aus Korollar 2.21 erhalten wir  $\hat{\beta}_{\text{ridge}} = (X^\top X + \lambda E_p)^{-1} (X^\top Y)$ , sodass

$$\begin{aligned} \hat{\beta}_{\text{ridge}} - \beta &= (X^\top X + \lambda E_p)^{-1} X^\top X\beta - \beta + (X^\top X + \lambda E_p)^{-1} X^\top \varepsilon \\ &= -(X^\top X + \lambda E_p)^{-1} \lambda \beta + (X^\top X + \lambda E_p)^{-1} X^\top \varepsilon. \end{aligned}$$



**Abb. 2.4** Quadrierte Residuen in Abhängigkeit von  $\beta = (\beta_1, \beta_2)^\top \in \mathbb{R}^2$  ohne Strafterm (*oben*) und mit  $\ell^2$ -Strafterm (*unten*)

Aus der Bias-Varianz-Zerlegung,  $\mathbb{E}[|Z|^2] = \text{tr}(\mathbb{E}[ZZ^\top])$  für beliebige Zufallsvektoren  $Z \in \mathbb{R}^p$  und dem Einsetzen von  $X^\top X = nE_p$  folgt

$$\begin{aligned} \mathbb{E}[|\hat{\beta}_{\text{ridge}} - \beta|^2] &= \lambda^2 |(X^\top X + \lambda E_p)^{-1} \beta|^2 \\ &\quad + \sigma^2 \text{tr}((X^\top X + \lambda E_p)^{-1} X^\top X (X^\top X + \lambda E_p)^{-1}) \\ &= \frac{|\beta|^2}{(1 + n\lambda^{-1})^2} + \frac{\sigma^2 p n}{(n + \lambda)^2}. \end{aligned}$$

Damit ergibt sich die behauptete Darstellung des mittleren quadratischen Fehlers von  $\hat{\beta}_{\text{ridge}}$ .  $\square$

Im Vergleich zum mittleren quadratischen Fehler  $\mathbb{E}[|\hat{\beta} - \beta|^2] = \sigma^2 p/n$  des Kleinst-Quadrate-Schätzers  $\hat{\beta}$  wird die Varianz des Ridge-Regressionsschätzers auf Kosten eines zusätzlichen Bias verringert. Da der Bias proportional zu  $|\beta|^2$  wächst, ist

der Shrinkageansatz insbesondere sinnvoll, wenn  $|\beta|$  klein ist. Die richtige Wahl des Penalisierungsparmeters  $\lambda$  ist hierbei allerdings entscheidend. Eine optimale Wahl von  $\lambda$  hängt vom unbekannten Vektor  $\beta$  ab und ist dem Statistiker nicht zugänglich (Aufgabe 2.9).

Der Ridge-Regressionschätzer liefert bei kleinen  $|\beta|^2$  auch noch gute Schätzergebnisse, wenn die Parameterdimension nicht viel kleiner als die Anzahl der Beobachtungen ist: Die in der Dimension linear wachsende Varianz  $\sigma^2 p$  des Kleinsten-Quadrate-Schätzers wird durch den Penalisierungsparmeter reduziert.

Die Ridge-Regression wird insbesondere bei schlecht konditionierter Designmatrix  $X$ , also wenn das Verhältnis von größtem zu kleinstem Eigenwert von  $X^\top X$  groß ist, verwendet. Dann ist einerseits die Lösung des Minimierungsproblems  $\hat{\beta}_{\text{ridge}} = (X^\top X + \lambda E_p)^{-1} X^\top Y$  numerisch stabiler, da die Matrix  $X^\top X + \lambda E_p$  besser konditioniert ist. Andererseits kann man analog zum Lemma beweisen, dass die Varianz des Schätzers insbesondere in Richtung der Eigenvektoren zu den kleinen Eigenwerten von  $X^\top X$  stark reduziert wird. Der Strafterm führt also zu einer numerischen wie statistischen Regularisierung.

Möchte man  $\beta$  in einem hochdimensionalen linearen Modell schätzen, also im Fall  $p \gg n$ , ist auch der Ridge-Regressionsschätzer nicht immer zielführend. Man kann aber den  $\ell^2$ -Strafterm durch einen  $\ell^1$ -Strafterm ersetzen, was typischerweise zu spärlich besetzten (englisch: *sparse*) Lösungen des Minimierungsproblems führt und auch bei sehr großen Parameterdimensionen  $p$  oft gut funktioniert. Die resultierende Schätzmethode ist der sogenannte Lasso-Schätzer, den wir in Kapitel ?? studieren werden.

### 2.1.3 Zufälliges Design und Vorhersage

Oft sind die Kovariablen im linearen Modell selbst zufällig und damit auch die Designmatrix  $X$ . Da wir neben der Zielgröße  $Y$  auch  $X$  beobachten, gelten die bisherigen Resultate zur Schätzung von  $\beta$  und  $\sigma$  weiter, sofern wir nur auf die Realisierung von  $X$  bedingen. Ein neuer wesentlicher Aspekt ergibt sich aber, wenn wir die Vorhersage der Zielgröße, gegeben die zugehörige Kovariable, betrachten.

**Definition 2.24** Im linearen Modell mit **zufälligem Design** beobachten wir  $(X_i, Y_i)_{i=1, \dots, n}$  mit

$$Y_i = \langle X_i, \beta \rangle + \varepsilon_i = (X\beta)_i + \varepsilon_i, \quad i = 1, \dots, n,$$

und der Designmatrix  $X = (X_1, \dots, X_n)^\top$  ( $X_i^\top$  ist die  $i$ -te Zeile in  $X$ ). Dabei seien  $\beta \in \mathbb{R}^p$  der unbekannte Parameter und  $(X_i, \varepsilon_i)_{i=1, \dots, n}$  unabhängige Zufallsvariablen mit identisch verteilten Zufallsvektoren  $X_i$  im  $\mathbb{R}^p$  und identisch verteilten reellwertigen Zufallsvariablen  $\varepsilon_i$ , für die  $\mathbb{E}[\varepsilon_i] = 0$ ,  $\sigma^2 := \mathbb{E}[\varepsilon_i^2] < \infty$  gilt und die unabhängig von  $X_i$  sind. Zusätzlich gelte  $\mathbb{E}[|X_i|^2] < \infty$ , sodass die **Design-Kovarianzmatrix**

$$\Sigma_X = \mathbb{E}[X_1 X_1^\top]$$

wohldefiniert ist.

Wegen der Unabhängigkeit der  $\varepsilon_i$  liegt ein gewöhnliches lineares Modell  $Y = X\beta + \varepsilon$  mit  $\Sigma = \sigma^2 E_n$  vor. Allerdings fordern wir hier nicht a priori, dass  $X$  vollen Rang  $p$  besitzt. Die Beobachtungen  $(X_i, Y_i)_{i=1, \dots, n}$  sind also i.i.d. mit  $\mathbb{E}[Y_i - \langle X_i, \beta \rangle] = 0$ . Man beachte, dass oft allgemeiner nur mit Hilfe der bedingten Erwartung  $\mathbb{E}[Y_i | X_i] = \langle X_i, \beta \rangle$ , also  $\mathbb{E}[\varepsilon_i | X_i] = 0$  gefordert wird, was bei uns aus der Unabhängigkeit von  $\varepsilon_i$  und  $X_i$  sowie  $\mathbb{E}[\varepsilon_i] = 0$  folgt. Falls  $X_1$  nicht zentriert ist, ist  $\Sigma_X$  nicht mehr die Kovarianzmatrix des Designs und man müsste eigentlich von der Matrix der zweiten Momente sprechen, was aber nicht gängig ist.

In der Sprache des maschinellen Lernens wollen wir nun auf der Grundlage der Trainingsdaten  $(X_i, Y_i)_{i=1, \dots, n}$  bei einer neuen Kovariablen (oder einem neuen Feature)  $X_{n+1}$  die Zielgröße  $Y_{n+1}$  vorhersagen. Im Folgenden schreiben wir  $\mathbb{E}^{(X_{n+1}, Y_{n+1})}$  für den Erwartungswert nur bezüglich der Zufallsvariablen  $X_{n+1}, Y_{n+1}$ .

**Definition 2.25** Für  $(X_i, Y_i)_{i=1, \dots, n+1}$  gelte das lineare Modell mit zufälligem Design. Für einen Schätzer  $\hat{\beta}_n$ , basierend auf  $(X_i, Y_i)_{i=1, \dots, n}$ , bezeichnet

$$\ell^{\text{pred}}(\hat{\beta}_n, \beta) = \mathbb{E}_{\beta}^{(X_{n+1}, Y_{n+1})} [(Y_{n+1} - \langle X_{n+1}, \hat{\beta}_n \rangle)^2]$$

den **Vorhersageverlust** (engl. *out-of-sample prediction loss*) und

$$R^{\text{pred}}(\hat{\beta}_n, \beta) = \mathbb{E}_{\beta}[\ell^{\text{pred}}(\hat{\beta}_n, \beta)] = \mathbb{E}_{\beta}[(Y_{n+1} - \langle X_{n+1}, \hat{\beta}_n \rangle)^2]$$

den **Vorhersagefehler** (auch Vorhersagerisiko, *out-of-sample prediction error*).

Der Vorhersageverlust ist also bezüglich der Realisierung von  $\hat{\beta}_n$  weiterhin zufällig, während der Vorhersagefehler als dessen Erwartungswert deterministisch ist. Im weiteren Verlauf werden wir einen kleinen Vorhersageverlust mit großer Wahrscheinlichkeit nachweisen können, während sich die Abschätzung seines Erwartungswerts als technisch schwierig herausstellt und wenig zusätzliche Aussagekraft besitzt.

**Lemma 2.26** Für den Vorhersageverlust gilt

$$\ell^{\text{pred}}(\hat{\beta}_n, \beta) = \langle \Sigma_X (\hat{\beta}_n - \beta), \hat{\beta}_n - \beta \rangle + \sigma^2.$$

**Beweis** Mit  $Y_{n+1} = \langle X_{n+1}, \beta \rangle + \varepsilon_{n+1}$  und  $\mathbb{E}[X_{n+1}\varepsilon_{n+1}] = \mathbb{E}[X_{n+1}]\mathbb{E}[\varepsilon_{n+1}] = 0$  folgt

$$\begin{aligned} \ell^{\text{pred}}(\hat{\beta}_n, \beta) &= \mathbb{E}^{(X_{n+1}, \varepsilon_{n+1})} [(\langle X_{n+1}, \beta - \hat{\beta}_n \rangle + \varepsilon_{n+1})^2] \\ &= (\beta - \hat{\beta}_n)^{\top} \mathbb{E}[X_{n+1} X_{n+1}^{\top}] (\beta - \hat{\beta}_n) + \mathbb{E}[\varepsilon_{n+1}^2] \\ &= (\beta - \hat{\beta}_n)^{\top} \Sigma_X (\beta - \hat{\beta}_n) + \sigma^2 \end{aligned}$$

und durch Umformulieren mit Skalarprodukt das Ergebnis.  $\square$

Da die Matrix  $\Sigma_X$  positiv semi-definit ist, ist der Vorhersageverlust immer mindestens  $\sigma^2$ , was für  $\hat{\beta}_n = \beta$  auch erreicht wird. Also selbst wenn wir  $\beta$  kennen,

machen wir noch einen Fehler bei der Vorhersage von  $Y_{n+1}$  auf Grund des stochastischen Fehlers in  $\varepsilon_{n+1}$ . Ein wichtiges Konzept des statistischen Lernens ist es, den nicht zugänglichen Vorhersagefehler, bei dem  $\beta$  und  $\Sigma_X$  nicht bekannt sind, durch einen empirischen Vorhersagefehler zu ersetzen und dann zu minimieren.

**Definition 2.27** Im linearen Modell mit zufälligem Design bezeichnet

$$R_n(\hat{\beta}_n) = \frac{1}{n} \sum_{i=1}^n (Y_i - \langle X_i, \hat{\beta}_n \rangle)^2$$

das **empirische Risiko** (engl. *empirical risk*, *in-sample risk*) eines Schätzers  $\hat{\beta}_n$ . Ein Schätzer  $\hat{\beta}_n$  mit

$$R_n(\hat{\beta}_n) = \inf_{\tilde{\beta}_n} R_n(\tilde{\beta}_n),$$

wobei das Infimum über alle Schätzer  $\tilde{\beta}_n$ , basierend auf  $(X_i, Y_i)_{i=1, \dots, n}$ , genommen wird, heißt **empirischer Risiko-Minimierer** (ERM, *empirical risk minimiser*).

Für deterministische  $b \in \mathbb{R}^p$  sehen wir sofort, dass  $\mathbb{E}_\beta[R_n(b)] = \ell^{\text{pred}}(b, \beta)$  gilt. Für den ERM  $\hat{\beta}_n$  statt  $b$  ist dies jedoch falsch.

**Lemma 2.28** *Kleinste Quadrate-Schätzer  $\hat{\beta}_n$  sind gerade die empirischen Risiko-Minimierer, und es gilt*

$$\mathbb{E}_\beta[R_n(\hat{\beta}_n)] = \frac{n-p}{n} \sigma^2,$$

sofern  $X$  fast sicher Rang  $p$  besitzt.

**Beweis** Beachte, dass  $R_n(\hat{\beta}_n) = \frac{1}{n} |Y - X\hat{\beta}_n|^2$  gilt. Die rechte Seite wird vom Kleinste-Quadrate-Schätzer minimiert, die linke Seite vom empirischen Risikominimierer, so dass beide Konzepte übereinstimmen. Aus Lemma 2.18 folgt dann

$$\mathbb{E}[R_n(\hat{\beta}_n)] = \frac{1}{n} \mathbb{E}[|Y - X\hat{\beta}_n|^2] = \frac{n-p}{n} \sigma^2.$$

Wegen der Unabhängigkeit von  $(X_i)$  und  $(\varepsilon_i)$  nimmt man dafür zunächst den Erwartungswert über  $(\varepsilon_i)$  gemäß Lemma 2.18 und erhält  $\frac{n-p}{n} \sigma^2$  fast sicher. Der Erwartungswert über  $(X_i)$  liefert dann den deterministischen Wert  $\frac{n-p}{n} \sigma^2$ .  $\square$

Eine einfache hinreichende Bedingung dafür, dass  $X$  vollen Rang  $p$  besitzt, ist, dass  $X_i$  gemäß einer Lebesguedichte im  $\mathbb{R}^p$  verteilt ist und  $n \geq p$  gilt. Das Argument beruht darauf, dass die Menge der  $p \times p$ -Matrizen  $A$  mit Determinante Null ( $A$  hat nicht vollen Rang) eine Lebesgue-Nullmenge im  $\mathbb{R}^{p \times p}$  ist.

Lemmata 2.26 und 2.28 zeigen sofort einen wichtigen Unterschied zwischen empirischem Risiko und Vorhersagefehler des empirischen Risikominimierers auf:

$$\mathbb{E}_\beta[R_n(\hat{\beta}_n)] = \frac{n-p}{n} \sigma^2 < \sigma^2 \leq R^{\text{pred}}(\hat{\beta}_n, \beta)$$

Da wir beim empirischen Risiko dieselben Daten für die Fehlerquantifizierung benutzen, die auch für den Schätzer verwendet wurden, ist der Fehler innerhalb der

Stichprobe stets kleiner als bei der Vorhersage außerhalb der Stichprobe. Wir wollen nun verstehen, inwiefern trotzdem die empirische Risikominimierung zu guten Schätzern bezüglich Vorhersagefehler führen, sofern der Stichprobenumfang  $n$  hinreichend groß ist. Man sagt auch, dass der Schätzer einen kleinen *Verallgemeinerungsfehler* (engl. *generalisation error*) besitzt, eine wichtige Frage auch später bei den Methoden des maschinellen Lernens.

**Satz 2.29** *Mit der empirischen Design-Kovarianzmatrix*

$$\Sigma_n = \frac{1}{n} \sum_{i=1}^n X_i X_i^\top = \frac{1}{n} X^\top X$$

gilt für den Vorhersageverlust und den Vorhersagefehler des Kleinste-Quadrate-Schätzers  $\widehat{\beta}_n$

$$\begin{aligned} \mathbb{E}^\varepsilon [\ell^{\text{pred}}(\widehat{\beta}_n, \beta)] &= \frac{\sigma^2}{n} \text{tr}(\Sigma_X \Sigma_n^{-1}) + \sigma^2, \\ R^{\text{pred}}(\widehat{\beta}_n, \beta) &= \frac{\sigma^2}{n} \mathbb{E}[\text{tr}(\Sigma_X \Sigma_n^{-1})] + \sigma^2, \end{aligned}$$

wobei  $\mathbb{E}^\varepsilon$  bedeutet, dass der Erwartungswert nur bezüglich  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$  genommen wird.

**Beweis** Es gilt  $\widehat{\beta}_n - \beta = (X^\top X)^{-1} X^\top \varepsilon = \Sigma_n^{-1} \frac{1}{n} X^\top \varepsilon$ , und wir erhalten mit Lemma 2.26 und  $\mathbb{E}[\varepsilon \varepsilon^\top] = \sigma^2 E_n$

$$\begin{aligned} \mathbb{E}^\varepsilon [\ell^{\text{pred}}(\widehat{\beta}_n, \beta)] &= \frac{1}{n^2} \mathbb{E}^\varepsilon [\langle \Sigma_X \Sigma_n^{-1} X^\top \varepsilon, \Sigma_n^{-1} X^\top \varepsilon \rangle] + \sigma^2 \\ &= \frac{1}{n^2} \text{tr}(\Sigma_X \Sigma_n^{-1} X^\top \mathbb{E}[\varepsilon \varepsilon^\top] X \Sigma_n^{-1}) + \sigma^2 \\ &= \frac{\sigma^2}{n} \text{tr}(\Sigma_X \Sigma_n^{-1} \Sigma_n \Sigma_n^{-1}) + \sigma^2, \end{aligned}$$

woraus die erste Behauptung folgt. Für den Vorhersagefehler verwenden wir dann  $R^{\text{pred}}(\widehat{\beta}_n, \beta) = \mathbb{E}^X [\mathbb{E}^\varepsilon [\ell^{\text{pred}}(\widehat{\beta}_n, \beta)]]$  wegen der Unabhängigkeit der  $X_i$  und  $\varepsilon_i$ .  $\square$

Im einfachen Fall eines Gaußschen Designs  $X_i \sim N(0, \Sigma_X)$  kann man mit Hilfe der sogenannten inversen Wishart-Verteilung zeigen, dass  $\mathbb{E}[\Sigma_n^{-1}] = \frac{n}{n-p-1} \Sigma_X^{-1}$  für  $n \geq p+2$  gilt, so dass  $R^{\text{pred}}(\widehat{\beta}_n, \beta) = \frac{\sigma^2 p}{n-p-1} + \sigma^2$  folgt. Im Allgemeinen ist nicht einmal klar, dass  $\Sigma_n^{-1}$  einen Erwartungswert besitzt. Die Theorie der zufälligen Matrizen liefert immerhin, dass  $\Sigma_n$  nahe an  $\Sigma_X$  ist, wenn  $n$  groß ist und große Werte der Kovariablen  $X_i$  ähnlich unwahrscheinlich sind wie unter einer Normalverteilung.

**Definition 2.30** Ein Zufallsvektor  $X$  im  $\mathbb{R}^p$  heißt  $(\Sigma, C)$ -**subgaußsch**, falls  $\mathbb{P}(|\langle X, u \rangle| > t) \leq C e^{-t^2/(2\langle \Sigma u, u \rangle)}$  für alle  $t > 0$ ,  $u \in \mathbb{R}^p \setminus \{0\}$  gilt mit einer Konstanten  $C > 0$  und einer positiv-definiten Matrix  $\Sigma \in \mathbb{R}^{p \times p}$ .

Insbesondere ist jeder normalverteilte Zufallsvektor  $X \sim N(\mu, \Sigma)$   $(\Sigma, C)$ -subgaußsch mit geeignetem  $C > 0$ , aber auch jeder Zufallsvektor, dessen Verteilung einen kompakten Träger besitzt. Ist mit der positiv-definiten Kovarianzmatrix  $\Sigma_X$

von  $X$  der standardisierte Zufallsvektor  $\Sigma_X^{-1/2} X$  ( $cE_p, C$ )-subgaußsch für ein  $c > 0$ , so ist  $X$  ( $c\Sigma_X, C$ )-subgaußsch, und umgekehrt. In diesem Sinn interpretieren wir im Folgenden  $\Sigma$  häufig als Vielfaches der Design-Kovarianzmatrix  $\Sigma_X$ . Ebenso wie normalverteilte Zufallsvektoren besitzen ( $c\Sigma_X, C$ )-subgaußsche Zufallsvektoren  $X$  endliche Momente, genauer gilt

$$\mathbb{E}[|\Sigma_X^{-1/2} X|^m]^{1/m} \leq \tilde{C}_m p^{1/2}, \quad \mathbb{E}[|\langle \Sigma_X^{-1/2} X, v \rangle|^m]^{1/m} \leq \tilde{C}_m |v|, \quad m \geq 1, \quad (2.4)$$

für deterministische  $v \in \mathbb{R}^p$ , vorausgesetzt  $\Sigma_X$  ist invertierbar. Die Konstanten  $\tilde{C}_m$  hängen nur von  $c, C$  und  $m$  ab.

Um die Abweichung der empirischen Design-Kovarianzmatrix  $\Sigma_n$  von der Design-Kovarianzmatrix  $\Sigma_X = \mathbb{E}[\Sigma_n]$  zu messen, stellt sich die *Spektralnorm* als geeignete Metrik heraus. Für eine Matrix  $M$  ist sie definiert als

$$\|M\| = \sup_{|v|=1} |Mv|.$$

Insbesondere gilt dann  $|Mu| \leq \|M\||u|$  für alle Vektoren  $u$ . Die Spektralnorm ist beschränkt durch die Frobeniusnorm, für die wir in Übung 2.11  $\mathbb{E}[\|\Sigma_n - \Sigma_X\|_2^2]^{1/2} = O(p/\sqrt{n})$  nachweisen. Ein erstaunliches Ergebnis der Theorie zufälliger Matrizen mit vielen Anwendungen in der hochdimensionalen Statistik ist, dass  $\|\Sigma_n - \Sigma_X\|$  unter milden Voraussetzungen von der stochastischen Ordnung  $\sqrt{p/n}$  ist. Insbesondere ist also  $\|\Sigma_n - \Sigma_X\|$  bereits klein, sobald die Dimension  $p$  von kleinerer Größenordnung als der Stichprobenumfang  $n$  ist. Wir passen Satz 4.7.1 von Vershynin (2018) auf unsere Situation an, siehe auch Satz 6.5 in Wainwright (2019).

**Satz 2.31 (Fehler der empirischen Kovarianzmatrix)** Sind  $X_1, \dots, X_n$  i.i.d. ( $c\Sigma_X, C$ )-subgaußsche Zufallsvektoren in  $\mathbb{R}^p$  mit  $C, c > 0$  und invertierbarer Kovarianzmatrix  $\Sigma_X \in \mathbb{R}^{p \times p}$ , so gilt für  $p \leq n$

$$\mathbb{E}[\|\Sigma_X^{-1/2} \Sigma_n \Sigma_X^{-1/2} - E_p\|] \leq \tilde{C} \sqrt{p/n}$$

mit einer Konstanten  $\tilde{C}$  (abhängig von  $c, C$ ).

Verwenden wir den  $O_{\mathbb{P}}$ -Kalkül aus Definition A.41, so gilt also

$$\|\Sigma_X^{-1/2} \Sigma_n \Sigma_X^{-1/2} - E_p\| = O_{\mathbb{P}}(\sqrt{p/n}).$$

Schreiben wir  $A := \Sigma_n^{1/2} \Sigma_X^{-1/2}$ , so erlaubt es dieses Resultat, weitere durch  $A$  definierte Matrizen abzuschätzen. Dazu nehmen wir  $\|A^T A - E_p\| \leq \varepsilon$  mit  $\varepsilon \in (0, 1)$  an. Die variationelle Charakterisierung von Eigenwerten symmetrischer Matrizen zeigt dann

$$\sup_{|v|=1} |Av|^2 - 1 = \sup_{|v|=1} |\langle (A^T A - E_p)v, v \rangle| = \|A^T A - E_p\| \leq \varepsilon.$$

Direkt folgt also  $\|A\| \leq \sqrt{1+\varepsilon} \leq 1 + \varepsilon/2$ . Darüberhinaus ist  $\inf_{|v|=1} |Av| \geq \sqrt{1-\varepsilon} > 0$ , so dass  $A$  invertierbar ist mit (setze  $w = A^{-1}v$ )

$$\|A^{-1}\| = \sup_{v \neq 0} \frac{|A^{-1}v|}{|v|} = \sup_{w \neq 0} \frac{|w|}{|Aw|} \leq \frac{1}{\sqrt{1-\varepsilon}} \leq 1 + \frac{\varepsilon}{1-\varepsilon}.$$

Mit denselben Argumenten schließen wir auf

$$\|(A^{-1})^\top A^{-1} - E_p\| = \sup_{|v|=1} ||A^{-1}v|^2 - 1| \leq \max\left(\frac{1}{1-\varepsilon} - 1, 1 - \frac{1}{1+\varepsilon}\right) = \frac{\varepsilon}{1-\varepsilon}.$$

Aus Satz 2.31 folgen daher durch Einsetzen und unter Verwendung von  $\|A^\top\| = \|A\|$  die weiteren asymptotischen Schranken unter der Asymptotik  $p/n \rightarrow 0$

$$\|\Sigma_n^{1/2} \Sigma_X^{-1/2}\| \leq 1 + O_{\mathbb{P}}(\sqrt{p/n}), \quad \|\Sigma_X^{1/2} \Sigma_n^{-1/2}\| \leq 1 + O_{\mathbb{P}}(\sqrt{p/n}), \quad (2.5)$$

$$\|\Sigma_n^{-1/2} \Sigma_X \Sigma_n^{-1/2} - E_p\| = O_{\mathbb{P}}(\sqrt{p/n}), \quad \|\Sigma_X^{1/2} \Sigma_n^{-1} \Sigma_X^{1/2} - E_p\| = O_{\mathbb{P}}(\sqrt{p/n}). \quad (2.6)$$

Insbesondere heißt das, dass  $\mathbb{P}(\Sigma_n^{-1} \text{ existiert}) \rightarrow 1$  gilt. Wir können nun immerhin den Vorhersageverlust des Kleinste-Quadrate-Schätzers in stochastischer Ordnung unter allgemeinen Bedingungen abschätzen. Dies ist für viele Zwecke vollkommen ausreichend und vermeidet zusätzliche Bedingungen, um die Existenz von  $\mathbb{E}[\Sigma_n^{-1}]$  sicherzustellen.

**Korollar 2.32** Die Kovariablen  $X_i$  seien  $(c\Sigma_X, C)$ -subgaußsche Zufallsvektoren im  $\mathbb{R}^p$  für Konstanten  $c, C > 0$ . Für den Vorhersageverlust des Kleinste-Quadrate-Schätzers  $\widehat{\beta}_n$  gilt

$$\ell^{\text{pred}}(\widehat{\beta}_n, \beta) = \frac{\sigma^2 p}{n} \left(1 + O_{\mathbb{P}}(\sqrt{p/n})\right) + \sigma^2,$$

sofern  $p \leq n$  und  $\Sigma_X$  invertierbar ist.

**Beweis** Aus Satz 2.29 folgt mit  $\text{tr}(E_p) = p$  und Kommutativität unter der Spur

$$\mathbb{E}^\varepsilon [\ell^{\text{pred}}(\widehat{\beta}_n, \beta)] = \frac{\sigma^2}{n} \left(p + \text{tr}(\Sigma_n^{-1/2} \Sigma_X \Sigma_n^{-1/2} - E_p)\right) + \sigma^2.$$

Wir verwenden  $|\text{tr}(M)| \leq p\|M\|$  für  $p \times p$ -Matrizen  $M$  sowie (2.6) und schließen

$$|\text{tr}(\Sigma_n^{-1/2} \Sigma_X \Sigma_n^{-1/2} - E_p)| = O_{\mathbb{P}}(p\sqrt{p/n}).$$

Also gilt

$$\mathbb{E}^\varepsilon [\ell^{\text{pred}}(\widehat{\beta}_n, \beta)] = \frac{\sigma^2 p}{n} \left(1 + O_{\mathbb{P}}(\sqrt{p/n})\right) + \sigma^2.$$

Da  $\ell^{\text{pred}}$  nicht-negativ ist, impliziert die Konvergenzrate des Erwartungswerts bezüglich  $\varepsilon$  dieselbe  $O_{\mathbb{P}}$ -Konvergenzrate und damit die Behauptung.  $\square$

Wir sehen, dass der zusätzliche Vorhersageverlust, neben dem unvermeidlichen Fehler  $\sigma^2$ , mit Rate  $\sigma^2 p/n$  gegen Null geht selbst bei wachsender Dimension  $p$ , falls nur  $p/n \rightarrow 0$  gilt.

Schließlich wollen wir noch den allgemeineren Zugang des statistischen Lernens kennenlernen. Der Hauptunterschied ist, dass die Daten mit einem sehr allgemeinen Modell beschrieben werden, für die aber nur eine Teilklasse von statistischen Methoden untersucht wird. Im vorliegenden Vorhersageproblem wird nurmehr angenommen, dass  $(X_i, Y_i)_{i=1, \dots, n+1}$  i.i.d. sind mit  $X_i \in \mathbb{R}^p$ ,  $Y_i \in \mathbb{R}$ . Wir betrachten die Klasse aller linearen Vorhersagen  $f_b(X_{n+1}) = \langle X_{n+1}, b \rangle$  von  $Y_{n+1}$  mit  $b \in \mathbb{R}^p$ . Der Vorhersageverlust ist dann

$$\ell^{\text{pred}}(b) = \mathbb{E}^{(X_{n+1}, Y_{n+1})} [(Y_{n+1} - f_b(X_{n+1}))^2],$$

was mit der vorherigen Definition von  $\ell^{\text{pred}}(b, \beta)$  übereinstimmt in dem Fall, dass wirklich das lineare Modell  $Y_{n+1} = \langle X_{n+1}, \beta \rangle + \varepsilon_{n+1}$  gilt. Im Allgemeinen gilt dies jedoch nicht, und wir vergleichen die Qualität eines Schätzers mit der besten Vorhersage innerhalb der linearen Klasse. Wir setzen daher

$$\beta^* := \arg \min_{b \in \mathbb{R}^p} \ell^{\text{pred}}(b)$$

und nennen  $\beta^*$  das *Orakel* oder den *Orakelparameter*. Da wir mit linearen Vorhersagen keinen geringeren Verlust als  $\ell^{\text{pred}}(f_{\beta^*})$  erzielen können, betrachten wir für einen Schätzer  $\hat{\beta}$  den *Exzessverlust*

$$\mathcal{E}^{\text{pred}}(\hat{\beta}) := \ell^{\text{pred}}(\hat{\beta}) - \ell^{\text{pred}}(\beta^*) = \ell^{\text{pred}}(\hat{\beta}) - \inf_{b \in \mathbb{R}^p} \ell^{\text{pred}}(b).$$

Mit dem empirischen Risiko von oben gilt weiterhin  $\mathbb{E}[R_n(b)] = \ell^{\text{pred}}(b)$  für deterministische  $b \in \mathbb{R}^p$ , und der Kleinste-Quadrate-Schätzer  $\hat{\beta}_n$  ist empirischer Risikominimierer, da die Berechnung ja unabhängig vom datenerzeugenden Modell ist. Ziel ist es nun, das Orakel  $\beta^*$  und die Größenordnung des Exzessrisikos von  $\hat{\beta}_n$  im allgemeinen Modell zu verstehen.

**Satz 2.33** Mit dem Kovarianzvektor  $C_{XY} := \mathbb{E}[X_1 Y_1]$  gilt für den Orakelparameter  $\beta^* = \Sigma_X^{-1} C_{XY}$ . Der Exzessverlust eines Schätzers  $\hat{\beta}$  ist gegeben durch

$$\mathcal{E}^{\text{pred}}(\hat{\beta}) = \langle \Sigma_X (\hat{\beta} - \beta^*), \hat{\beta} - \beta^* \rangle, \quad (2.7)$$

wobei  $\Sigma_X^{-1}$  existieren möge. Sind ferner die Kovariablen  $X_i$   $(c\Sigma_X, C)$ -subgaußsch, gilt  $\mathbb{E}[Y_1^4] < \infty$  und ist  $n \geq p$ , so gilt für den Kleinste-Quadrate-Schätzer  $\hat{\beta}_n$

$$\mathcal{E}^{\text{pred}}(\hat{\beta}_n) = O_{\mathbb{P}}(p/n),$$

wobei die Konstante nur von  $c$ ,  $C$  und  $\mathbb{E}[Y_1^4]$  abhängt.

**Beweis** Durch Einsetzen und quadratische Ergänzung erhalten wir

$$\ell^{\text{pred}}(b) = \mathbb{E}[Y_{n+1}^2] - 2\mathbb{E}[Y_{n+1} \langle X_{n+1}, b \rangle] + \mathbb{E}[\langle X_{n+1}, b \rangle^2]$$

$$\begin{aligned}
&= \mathbb{E}[Y_1^2] - \langle C_{XY}, b \rangle + \langle \Sigma_X b, b \rangle \\
&= \langle \Sigma_X(b - \Sigma_X^{-1} C_{XY}), b - \Sigma_X^{-1} C_{XY} \rangle + \mathbb{E}[Y_1^2] - \langle \Sigma_X^{-1} C_{XY}, C_{XY} \rangle.
\end{aligned}$$

Dies wird offensichtlich durch  $\beta^* = \Sigma_X^{-1} C_{XY}$  in  $b$  minimiert. Weiterhin folgt aus der Darstellung sofort (2.7) für den Exzessverlust.

Zur asymptotischen Analyse von  $\hat{\beta}_n = \Sigma_n^{-1} C_n$  mit

$$C_n = \frac{1}{n} X^\top Y = \frac{1}{n} \sum_{i=1}^n X_i Y_i$$

wollen wir zunächst Erwartungswert und Varianz von  $\Sigma_X^{-1/2} C_n$  abschätzen. Es gilt  $\mathbb{E}[C_n] = C_{XY}$  sowie mit der Cauchy-Schwarz-Ungleichung

$$\begin{aligned}
|\Sigma_X^{-1/2} C_{XY}|^2 &= \sup_{|v|=1} \langle \Sigma_X^{-1/2} X_1 Y_1, v \rangle^2 \\
&= \sup_{|v|=1} \mathbb{E}[\langle \Sigma_X^{-1/2} X_1, v \rangle Y_1]^2 \\
&\leq \sup_{|v|=1} \mathbb{E}[\langle \Sigma_X^{-1/2} X_1, v \rangle^2] \mathbb{E}[Y_1^2] \\
&\leq \tilde{C}_2^2 \mathbb{E}[Y_1^2]
\end{aligned}$$

mit  $\tilde{C}_2$  aus (2.4). Wegen  $X_i Y_i$  i.i.d. folgt ferner für die Varianz des Zufallsvektors  $\Sigma_X^{-1/2} C_n$  aus der Cauchy-Schwarz-Ungleichung

$$\begin{aligned}
\text{Var}(\Sigma_X^{-1/2} C_n) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\Sigma_X^{-1/2} X_i Y_i) \leq \frac{1}{n} \mathbb{E}[|\Sigma_X^{-1/2} X_1 Y_1|^2] \\
&\leq \frac{1}{n} \mathbb{E}[|\Sigma_X^{-1/2} X_1|^4]^{1/2} \mathbb{E}[Y_1^4]^{1/2}.
\end{aligned}$$

Unter Verwendung von  $\tilde{C}_4$  aus (2.4) und  $\mathbb{E}[Y_1^2]^2 \leq \mathbb{E}[Y_1^4]$  implizieren beide Abschätzungen die stochastischen Ordnungen

$$|\Sigma_X^{-1/2} (C_n - C_{XY})| = O_{\mathbb{P}}(\sqrt{p/n}), \quad |\Sigma_X^{-1/2} C_n| = O_{\mathbb{P}}(1)$$

mit einem Faktor in Abhängigkeit von  $c, C, \mathbb{E}[Y_1^4]$ . Einsetzen in (2.7) und Verwenden von (2.6) ergibt damit

$$\begin{aligned}
\mathcal{E}^{\text{pred}}(\hat{\beta}_n) &= |\Sigma_X^{1/2}(\hat{\beta}_n - \beta^*)|^2 \\
&= |\Sigma_X^{1/2} \Sigma_n^{-1} C_n - \Sigma_X^{-1/2} C_{XY}|^2 \\
&\leq \left( \|\Sigma_X^{1/2} \Sigma_n^{-1} \Sigma_X^{1/2} - E_p\| |\Sigma_X^{-1/2} C_n| + |\Sigma_X^{-1/2} (C_n - C_{XY})| \right)^2 \\
&= \left( O_{\mathbb{P}}(\sqrt{p/n}) O_{\mathbb{P}}(1) + O_{\mathbb{P}}(\sqrt{p/n}) \right)^2 = O_{\mathbb{P}}(p/n),
\end{aligned}$$

wie behauptet.  $\square$

Wenn wir mit Korollar 2.32 im korrekt spezifizierten linearen Modell vergleichen, so sehen wir, dass in beiden Fällen die Rate  $p/n$  für den Exzessverlust vorliegt (beachte  $\inf_b \ell^{\text{pred}}(b, \beta) = \sigma^2$ ). Allerdings ist dort  $\sigma^2 p/n$  die exakte Form des Terms größter Ordnung, während der Faktor im allgemeinen Fall unbestimmt bleibt. Trotzdem garantiert dieses Resultat eine Robustheit gegenüber Modellmisspezifikation, also falls die Annahme eines linearen Modellzusammenhangs verletzt ist.

## 2.2 Inferenz unter Normalverteilungsannahme

Statistische Inferenz umfasst die Konstruktion von Tests und Konfidenzintervallen. Im Gegensatz zur Schätztheorie aus dem letzten Abschnitt, benötigen wir hier eine explizite Annahme an die Verteilung des Fehlervektors  $\varepsilon$  im linearen Modell. Im Folgenden werden wir stets das gewöhnliche lineare Modell unter der Normalverteilungsannahme ( $\varepsilon_i \sim N(0, \sigma^2 E_n)$ ) betrachten, da die Abweichungen der (Mess-)Werte vieler natur-, wirtschafts- und ingenieurwissenschaftlicher Vorgänge durch die Normalverteilung in guter Näherung beschrieben werden können. Zusätzlich beschränken wir uns auf den Fall, dass die Designmatrix  $X$  deterministisch ist und vom Rang  $p < n$ .

Ist die Fehlervarianz  $\sigma^2$  bekannt, lassen sich unter der Normalverteilungsannahme leicht Konfidenzintervalle konstruieren:

*Beispiel 2.34* Sind die Messfehler ( $\varepsilon_i \sim N(0, \sigma^2 E_n)$ ) gemeinsam normalverteilt und  $\rho = \langle v, \beta \rangle$  für  $v \in \mathbb{R}^k$ , so gilt

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^\top X)^{-1}) \quad \text{und} \quad \hat{\rho} = \langle v, \hat{\beta} \rangle \sim N(\rho, \sigma^2 \langle (X^\top X)^{-1} v, v \rangle).$$

Ist  $\sigma > 0$  bekannt, so ist ein Konfidenzintervall zum Niveau  $1 - \alpha$  für  $\rho$  gegeben durch

$$[\hat{\rho} - q_{1-\alpha/2} \sigma \sqrt{\langle (X^\top X)^{-1} v, v \rangle}, \hat{\rho} + q_{1-\alpha/2} \sigma \sqrt{\langle (X^\top X)^{-1} v, v \rangle}],$$

mit dem  $(1 - \alpha/2)$ -Quantil  $q_{1-\alpha/2}$  der Standardnormalverteilung. Beachte, dass dieses Konfidenzintervall über den Korrespondenzsatz 1.70 dem zweiseitigen Gauß-Test aus Beispiel 1.65 entspricht. Insbesondere ergibt sich für  $\alpha = 0,05$  der Wert  $q_{0,975} \approx 1,96$ , der in Anwendungen häufig auftaucht.

Die Annahme in diesem Beispiel, dass  $\sigma$  bekannt sei, ist in den wenigstens Fällen erfüllt. Bei unbekanntem Rauschniveau können wir  $\sigma$  durch den Schätzer  $\hat{\sigma}$  ersetzen. Ist  $\hat{\sigma}$  konsistent, dann folgt aus Slutzkys Lemma, dass das resultierende Konfidenzintervall beziehungsweise der entsprechende Test das vorgegebene Niveau zumindest asymptotisch erreicht. Im nächsten Abschnitt werden wir sehen, dass man sogar die Verteilung der normalisierten Statistiken explizit bestimmen und so Konfidenzbereiche und Tests konstruieren kann, die auch für endliche Stichprobengrößen ein gegebenes Niveau genau erreichen.

Folgende Verteilungen sind essentiell für die angestrebte Inferenz und bilden die Grundlage für die vielfach genutzten  $t$ - und  $F$ -Tests.

**Definition 2.35** Die **t-Verteilung**  $t(n)$  (oder Student-t-Verteilung) mit  $n \in \mathbb{N}$  Freiheitsgraden auf  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  ist gegeben durch die Lebesgue-Dichte

$$t_n(x) = \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})\sqrt{\pi n}} \left(1 + \frac{x^2}{n}\right)^{-(n+1)/2}, \quad x \in \mathbb{R}.$$

Dabei bezeichnet  $\Gamma(p) = \int_0^\infty t^{p-1} e^{-t} dt$  die Gammafunktion.

**Definition 2.36** Die **F-Verteilung**  $F(m, n)$  (oder Fisher-Verteilung) mit  $(m, n) \in \mathbb{N}^2$  Freiheitsgraden auf  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  ist gegeben durch die Lebesgue-Dichte

$$f_{m,n}(x) = \frac{m^{m/2} n^{n/2}}{B(\frac{m}{2}, \frac{n}{2})} \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}} \mathbb{1}_{\mathbb{R}^+}(x), \quad x \in \mathbb{R}.$$

Dabei bezeichnet  $B(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$  die Betafunktion.

**Lemma 2.37** Es seien  $X_1, \dots, X_m, Y_1, \dots, Y_n$  unabhängige  $N(0, 1)$ -verteilte Zufallsvariablen. Dann gilt

$$T_n := \frac{X_1}{\sqrt{\frac{1}{n} \sum_{j=1}^n Y_j^2}} \sim t(n) \quad \text{und} \quad F_{m,n} := \frac{\frac{1}{m} \sum_{i=1}^m X_i^2}{\frac{1}{n} \sum_{j=1}^n Y_j^2} \sim F(m, n).$$

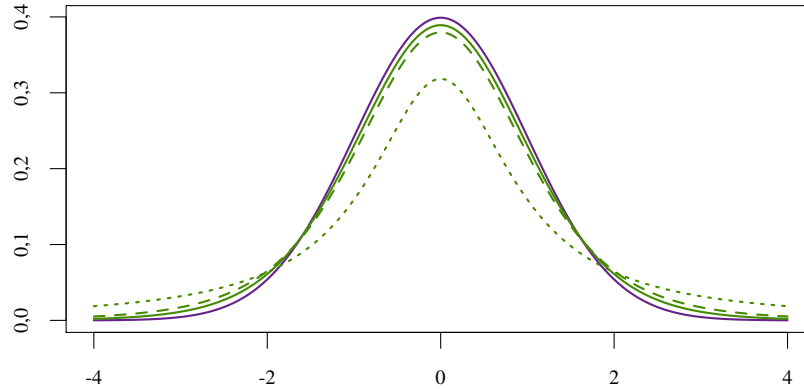
**Beweis** Es gilt  $T_n^2 = F_{1,n}$ , sodass mittels Dichtetransformation  $f^{|T_n|}(x) = f^{F_{1,n}}(x^2)2x$ ,  $x \geq 0$ , folgt. Da  $T_n$  symmetrisch (wie  $-T_n$ ) verteilt ist, folgt  $f^{T_n}(x) = f^{F_{1,n}}(x^2)|x|$ ,  $x \in \mathbb{R}$ , und Einsetzen zeigt die Behauptung für  $T_n$ , sofern  $F_{1,n}$   $F(1, n)$ -verteilt ist.

Um die Behauptung für  $F_{m,n}$  nachzuweisen, benutzen wir, dass  $X := \sum_{i=1}^m X_i^2$  bzw.  $Y := \sum_{j=1}^n Y_j^2$  gemäß  $\chi^2(m)$  bzw.  $\chi^2(n)$  verteilt sind. Wegen der Unabhängigkeit von  $X$  und  $Y$  und des Satzes von Fubini gilt für  $z > 0$  (setze  $w = x/y$ )

$$\begin{aligned} \mathbb{P}(X/Y \leq z) &= \int_0^\infty \int_0^\infty \mathbb{1}_{\{x/y \leq z\}} f^X(x) f^Y(y) dx dy \\ &= \int_0^z \left( \int_0^\infty f^X(wy) f^Y(y) y dy \right) dw. \end{aligned}$$

Setzen wir die  $\chi^2$ -Dichten ein und substituieren  $w = (z+1)y$  und  $t = w/2$ , ergibt sich die Dichte von  $X/Y$ :

$$\begin{aligned} f^{X/Y}(z) &= \int_0^\infty f^X(zy) f^Y(y) y dy \\ &= \frac{2^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty (zy)^{m/2-1} y^{n/2} e^{-(zy+y)/2} dy \end{aligned}$$



**Abb. 2.5** Für  $n \rightarrow \infty$  konvergiert die Dichte der  $t(n)$ -Verteilung (grün) gegen jene der Standardnormalverteilung (violett). Gepunktet entspricht  $t(1)$ , gestrichelt  $t(5)$  und durchgezogen  $t(10)$ .

$$\begin{aligned}
 &= \frac{2^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty (zw/(z+1))^{m/2-1} (w/(z+1))^{n/2} e^{-w/2} (z+1)^{-1} dw \\
 &= \frac{z^{m/2-1} (z+1)^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty 2^{-((m+n)/2-1)} 2^{-1} w^{(m+n)/2-1} e^{-w/2} dw \\
 &= \frac{z^{m/2-1} (z+1)^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty t^{((m+n)/2-1)} e^{-t} dt \\
 &= \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} z^{m/2-1} (z+1)^{-(m+n)/2}, \quad z > 0.
 \end{aligned}$$

Dichtetransformation ergibt damit für  $F_{m,n} = \frac{n}{m} \frac{X}{Y}$  die Dichte  $\frac{m}{n} f^{X/Y}(\frac{m}{n}x) = f_{m,n}(x)$ .  $\square$

**Bemerkung 2.38 ( $t$ - und  $F$ -Verteilung)** Die  $t$ -Verteilung ist symmetrisch zur Symmetrieachse  $x = 0$  und glockenförmig. Die Dichte fällt jedoch nur polynomiell statt exponentiell schnell ab. Insbesondere besitzt die  $t(n)$ -Verteilung für jedes  $n \in \mathbb{N}$  nur Momente bis zur Ordnung  $p < n$ . Man spricht von *heavy tails*, was man mit „schweren Flanken“ übersetzen könnte. Für  $n = 1$  ist die  $t(n)$ -Verteilung gerade die Cauchy-Verteilung, und für  $n \rightarrow \infty$  konvergiert sie schwach gegen die Standardnormalverteilung, was man leicht aus dem letzten Lemma folgern kann, siehe auch Abbildung 2.5. Die  $F$ -Verteilung hat ebenfalls heavy tails. Es gilt  $F_{1,n} = T_n^2$ . Für  $n \rightarrow \infty$  konvergiert  $mF_{m,n}$  gegen die  $\chi^2(m)$ -Verteilung.

Folgendes Hilfsresultat zur Verteilung von quadratischen Formen wird uns auch bei der Konstruktion von Tests und Konfidenzbändern im linearen Modell helfen:

**Lemma 2.39** Sei  $Z \sim N(0, E_n)$  und sei  $R \in \mathbb{R}^{n \times n}$  eine Orthogonalprojektion vom Rang  $\text{rank}(R) = r \leq n$ . Dann gilt

- (i)  $Z^\top RZ \sim \chi^2(r)$ ,  
(ii)  $Z^\top RZ$  ist unabhängig von  $BZ$  für jede Matrix  $B \in \mathbb{R}^{p \times n}$  mit  $BR = 0$ ,  
(iii) für jede weitere Orthogonalprojektion  $S \in \mathbb{R}^{n \times n}$  mit  $\text{rank}(S) = s \leq n$  und  $RS = 0$  sind  $Z^\top RZ$  und  $Z^\top SZ$  unabhängig, und es gilt

$$\frac{s}{r} \frac{Z^\top RZ}{Z^\top SZ} \sim F(r, s).$$

**Beweis** (i) Als Orthogonalprojektion ist  $R$  symmetrisch und idempotent ( $R = R^\top$  und  $R^2 = R$ ). Daher existiert eine Orthogonalmatrix  $T$  mit

$$R = TD_r T^\top \text{ und } D_r = \begin{pmatrix} E_r & 0 \\ 0 & 0 \end{pmatrix}.$$

Da  $T$  orthogonal ist und  $Z$  standardnormalverteilt, folgt  $W := T^\top Z \sim N(0, E_n)$ . Wegen

$$Z^\top RZ = Z^\top (TD_r T^\top)Z = (T^\top Z)^\top D_r (T^\top Z) = W^\top D_r W = \sum_{i=1}^r W_i^2$$

ist  $Z^\top RZ \sim \chi^2(r)$ -verteilt.

(ii) Wir setzen  $Y := BZ \sim N(0, B^\top B)$  und  $V := RZ \sim N(0, R)$ . Dann gilt

$$\text{Cov}(Y, V) = B \text{Cov}(Z) R^\top = BR = 0.$$

Weiter haben wir  $Z^\top RZ = Z^\top R^2 Z = (RZ)^\top (RZ) = V^\top V$ . Da  $(Y, V)$  als Lineartransformation von  $Z$  gemeinsam normalverteilt ist, folgt aus der Unkorreliertheit bereits die Unabhängigkeit von  $Y$  und  $V$  und somit auch die von  $Y = BZ$  und  $V^\top V = Z^\top RZ$ .

(iii) Genau wie in (ii) folgt die Unabhängigkeit von  $Y := SZ$  und  $V := RZ$  und somit auch die Unabhängigkeit von  $Y^\top Y = Z^\top SZ$  und  $V^\top V = Z^\top RZ$ . Zusammen mit (i) und dem vorangegangenen Lemma folgt die Behauptung.  $\square$

Als Folgerung erhalten wir Tests und Konfidenzbereiche für die Schätzung von  $\beta$  und linearen Funktionalen im gewöhnlichen linearen Modell unter der Normalverteilungsannahme:

**Satz 2.40 (F-Test)** Betrachte im gewöhnlichen linearen Modell unter der Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  mit unbekanntem  $\sigma > 0$  das Testproblem  $H_0: \beta = \beta_0, \sigma > 0$  gegen  $H_1: \beta \neq \beta_0, \sigma > 0$ . Dann ist der zweiseitige **F-Test** (auch **Fisher-Test**)

$$\varphi_\alpha = \mathbb{1}(F > q_{F(p, n-p); 1-\alpha}) \quad \text{mit} \quad F = \frac{\frac{1}{p} |X(\hat{\beta} - \beta_0)|^2}{\hat{\sigma}^2},$$

dem Kleinst-Quadrate-Schätzer  $\widehat{\beta}$ , der empirischen Stichprobenvarianz  $\widehat{\sigma}^2 = \frac{1}{n-p} |Y - X\widehat{\beta}|^2$  und dem  $(1 - \alpha)$ -Quantil  $q_{F(p, n-p); 1-\alpha}$  der  $F(p, n-p)$ -Verteilung ein Likelihood-Quotiententest zum Niveau  $\alpha \in (0, 1)$ .

**Beweis** Da der Kleinst-Quadrate-Schätzer gerade der Maximum-Likelihood-Schätzer unter Normalverteilung ist, können wir die Likelihood für alle  $\sigma^2 > 0$  über

$$\sup_{\beta \in \mathbb{R}^p \setminus \{\beta_0\}} L(\sigma^2, \beta) = L(\sigma^2, \widehat{\beta})$$

maximieren. Wir erhalten für den Likelihood-Quotienten (Methode 1.63)

$$T = \frac{\sup_{\sigma^2 > 0} L(\widehat{\beta}, \sigma^2)}{\sup_{\sigma^2 > 0} L(\beta_0, \sigma^2)} = \frac{\sup_{\sigma^2 > 0} \sigma^{-n} \exp(-\frac{1}{2\sigma^2} |Y - X\widehat{\beta}|^2)}{\sup_{\sigma^2 > 0} \sigma^{-n} \exp(-\frac{1}{2\sigma^2} |Y - X\beta_0|^2)}.$$

Wegen

$$\frac{d}{d\sigma^2} \left( -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} |Y - X\widehat{\beta}|^2 \right) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} |Y - X\widehat{\beta}|^2,$$

wird das Supremum im Zähler von  $T$  bei  $\widehat{\sigma}_{\text{MLE},1}^2 = \frac{1}{n} |Y - X\widehat{\beta}|^2$  angenommen. Analog wird das Supremum im Nenner von  $T$  bei  $\widehat{\sigma}_{\text{MLE},0}^2 = \frac{1}{n} |Y - X\beta_0|^2$  angenommen. Mit  $X\widehat{\beta} = \Pi_X Y$  und der Orthogonalität  $(E_n - \Pi_X)\Pi_X = 0$  erhalten wir

$$\begin{aligned} T &= \frac{\widehat{\sigma}_{\text{MLE},1}^{-n}}{\widehat{\sigma}_{\text{MLE},0}^{-n}} = \left( \frac{|Y - X\beta_0|^2}{|Y - \Pi_X Y|^2} \right)^{n/2} \\ &= \left( \frac{|(E_n - \Pi_X)Y|^2 + |\Pi_X(Y - X\beta_0)|^2}{|Y - \Pi_X Y|^2} \right)^{n/2} = \left( 1 + \frac{|X(\widehat{\beta} - \beta_0)|^2}{(n-p)\widehat{\sigma}^2} \right)^{n/2} \end{aligned}$$

mit der empirischen Stichprobenvarianz  $\widehat{\sigma}^2 = |Y - \Pi_X Y|^2 / (n-p)$ . Durch monotone Transformation hat also ein Likelihood-Quotiententest die Form

$$\varphi_\alpha = \mathbb{1} \left( \frac{\frac{1}{p} |X(\widehat{\beta} - \beta_0)|^2}{\widehat{\sigma}^2} > \widetilde{c}_\alpha \right),$$

wobei wir wegen der stetigen Verteilung auf die Randomisierung verzichten können. Da  $\Pi_X$  und  $E_n - \Pi_X$  als Projektionen auf  $\text{Im}(X)$  bzw.  $(\text{Im}(X))^\perp$  symmetrische, idempotente Matrizen mit Rang  $p$  bzw.  $(n-p)$  sind und  $(E_n - \Pi_X)\Pi_X = 0$  gilt, folgt aus Lemma 2.39 unter  $H_0$ :

$$\frac{\frac{1}{p} |X(\widehat{\beta} - \beta_0)|^2}{\widehat{\sigma}^2} = \frac{(n-p)}{p} \frac{\varepsilon^\top \Pi_X \varepsilon}{\varepsilon^\top (E_n - \Pi_X) \varepsilon} \sim F(p, n-p).$$

Man beachte, dass sich der wahre Wert von  $\sigma^2$  im Bruch herauskürzt. Wählen wir also  $\widetilde{c}_\alpha = q_{F(p, n-p); 1-\alpha}$ , so besitzt  $\varphi_\alpha$  Niveau  $\alpha$  unter  $H_0$ .  $\square$

Im Allgemeinen bezeichnen wir einen Hypothesentest als F-Test, wenn seine Teststatistik unter der Nullhypothese einer F-Verteilung folgt. Zwei Hauptanwendungen von F-Tests sind die Überprüfung, ob die Regressionskoeffizienten in der linearen Regression signifikant von vorgegebenen Koeffizienten abweichen, und der Nachweis, ob sich die Mittelwerte aus zwei oder mehr Stichproben aus unterschiedlichen, normalverteilten Populationen signifikant unterscheiden (Varianzanalyse, Methode 2.62).

*Beispiel 2.41 (F-Test, Happiness-Score)* Erinnern wir uns an den World Happiness Report und speziell an den Zusammenhang zwischen pro-Kopf-Bruttoinlandsprodukt und Happiness-Score (siehe Beispiel 1.3). Ein Glücksforscher gibt als Faustregel die Formel  $y = 2x + 4$  für den Happiness-Score  $y$  und das pro-Kopf-Bruttoinlandsprodukt  $x$  an. Mit  $p = 2$  und  $n = 156$  ist die Hypothese damit  $H_0: \beta = \beta_0$  für  $\beta_0 = (4, 2)^\top$ . Wir legen das Signifikanzniveau auf  $\alpha = 0,05$  fest. Unter Zuhilfenahme einer Statistik-Software bestimmen wir den Kleinste-Quadrate-Schätzer, die Stichprobenvarianz, das entsprechende  $(1 - \alpha)$ -Quantil der  $F_{2,198}$ -Verteilung sowie den Wert der F-Statistik:

$$\hat{\beta} \approx \begin{pmatrix} 3,40 \\ 2,22 \end{pmatrix}, \quad \hat{\sigma}^2 \approx 0,46, \quad q_{F(p,n-p);1-\alpha} \approx 3,05 \quad \text{und} \quad F \approx 28,77$$

Wegen  $F > q_{F(p,n-p);1-\alpha}$  wird die Nullhypothese verworfen. Die Faustregel ist demnach eine zu starke Vereinfachung der tatsächlichen Parameterwerte.

Für das Testen des reellen Parameters  $\rho = \langle v, \beta \rangle$  zu gegebenem  $v \in \mathbb{R}^p$  erhalten wir:

**Satz 2.42 (t-Test)** *Im gewöhnlichen linearen Modell unter Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  mit unbekanntem  $\sigma > 0$  betrachten wir den abgeleiteten Parameter  $\rho = \langle v, \beta \rangle$  für ein  $v \in \mathbb{R}^p \setminus \{0\}$ . Dann ist der Likelihood-Quotiententest der Hypothese  $H_0: \rho = \rho_0, \sigma > 0$  gegen die Alternative  $H_1: \rho \neq \rho_0, \sigma > 0$  für ein  $\rho_0 \in \mathbb{R}$  zum Niveau  $\alpha \in (0, 1)$  gegeben durch den zweiseitigen **t-Test** (auch **Student-t-Test**)*

$$\varphi_\alpha = \mathbb{1}(|T| > q_{t(n-p);1-\alpha/2}) \quad \text{mit} \quad T := \frac{\hat{\rho} - \rho_0}{\hat{\sigma} \sqrt{\langle (X^\top X)^{-1} v, v \rangle}}$$

und  $\hat{\rho} = \langle v, \hat{\beta} \rangle$ , dem Kleinste-Quadrate-Schätzer  $\hat{\beta}$ , der empirischen Stichprobenvarianz  $\hat{\sigma}^2$  und dem  $(1 - \alpha/2)$ -Quantil  $q_{t(n-p);1-\alpha/2}$  der  $t(n-p)$ -Verteilung.

Im Allgemeinen bezeichnen wir Hypothesentests, deren Teststatistik unter der Nullhypothese t-verteilt sind, als t-Tests. Neben dem Einstichproben-t-Test gibt es auch den *Zweistichproben-t-Test*, der die Mittelwerte zweier unabhängiger Stichproben auf Gleichheit testet (siehe Korollar 2.65).

**Beweis** Der direkte Nachweis, dass der angegebene t-Test ein Likelihood-Quotiententest ist, verbleibt als Übung, siehe Aufgabe 2.13. Später wird diese Aussage in Beispiel 2.51 auch aus der allgemeinen Theorie folgen.

Aus dem Satz von Gauß-Markov und der Normalverteilungsannahme folgt  $\hat{\rho} \sim N(\rho_0, \sigma^2 v^\top (X^\top X)^{-1} v)$  unter  $H_0: (\beta, \sigma) \in \{b \in \mathbb{R}^p : \langle v, b \rangle = \rho_0\} \times (0, \infty)$ . Daraus folgt

$$\frac{\hat{\rho} - \rho_0}{\sigma \sqrt{\langle (X^\top X)^{-1} v, v \rangle}} \sim N(0, 1) \text{ unter } H_0.$$

Andererseits sind  $\hat{\rho}$  und  $\hat{\sigma}^2$  unabhängig (wegen  $(E_n - \Pi_X)\Pi_X = 0$  und Lemma 2.39), und es gilt  $\hat{\sigma}^2 = \sigma^2 Y / (n - p)$  für eine Zufallsvariable  $Y \sim \chi^2(n - p)$ . Damit ist

$$\frac{\hat{\rho} - \rho_0}{\sqrt{\hat{\sigma}^2 \langle (X^\top X)^{-1} v, v \rangle}} \sim t(n - p) \text{ unter } H_0,$$

und die Behauptung folgt durch Wahl des richtigen Quantils.  $\square$

**Beispiel 2.43 (Einstichproben-t-Test)** Mit dem Einstichproben-t-Test wird auf den Mittelwert einer i.i.d. normalverteilten Stichprobe getestet. Mögliche Anwendungsfälle sind Nachweise, ob Sollwerte, zum Beispiel das Füllgewicht von Zuckerpackungen, eingehalten werden oder die Istwerte signifikant davon abweichen.

Wir beobachten also  $Y_i = \beta + \varepsilon_i$  mit  $\varepsilon_i \sim N(0, \sigma^2)$  und unbekanntem Mittelwert  $\beta$ . Wir erhalten ein lineares Modell mit  $X = (1, \dots, 1)^\top \in \mathbb{R}^n$  und dem empirischen Mittelwert als Schätzer für  $\beta$ :

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y = \frac{1}{n} \sum_{i=1}^n Y_i$$

Wir wenden nun Satz 2.42 an, wobei  $p = 1$  und  $v = 1$ , sodass  $\rho = \beta$  und  $\hat{\rho} = \hat{\beta}$ . Wir erhalten die Teststatistik

$$T = \frac{\sqrt{n}(\bar{Y} - \rho_0)}{\hat{\sigma}}.$$

Im Gegensatz zum Gauß-Test aus Beispiel 2.34 wird also im t-Test die Varianz  $\sigma^2$  durch die Stichprobenvarianz  $\hat{\sigma}^2$  ersetzt und ein Quantil der  $t$ -Verteilung benutzt.

**Kurzbiografie (William Sealy Gosset)** William Sealy Gosset wurde 1876 in Canterbury (Südostengland) geboren und studierte Chemie und Mathematik am New College in Oxford. Im Anschluss begann er bei der Dubliner Brauerei Arthur Guinness & Son zu arbeiten. Gossets Augenmerk war auf statistische Tests mit kleinen Stichprobengrößen gerichtet, was ein typisches Problem von Brauereien war. Da die Guinness-Brauerei ihren Mitarbeitern verbot, Arbeiten zu veröffentlichen, publizierte Gosset seine Erkenntnisse unter dem Pseudonym *Student* – daher auch der Name *Student-t-Verteilung*. Insbesondere Ronald Aylmer Fisher erkannte die Bedeutung von Gossets Arbeiten und entwickelte sie weiter, woraus die *t-Teststatistik* und die Anwendung der  $t$ -Verteilung in der Regressionsanalyse entstand. 1935 nahm Gosset eine Führungsposition in der neuen Guinness-Brauerei in London an, starb jedoch schon zwei Jahre später.

Der Korrespondenzsatz 1.70 zwischen Tests und Konfidenzbereichen liefert uns sofort folgende wichtige Konstruktionen von Konfidenzmengen:

**Satz 2.44 (Konfidenzmengen im linearen Modell)** Im gewöhnlichen linearen Modell unter der Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  für  $\sigma > 0$  gelten für den Kleinste-Quadrate-Schätzer  $\hat{\beta}$  und die empirische Stichprobenvarianz  $\hat{\sigma}^2$  folgende Konfidenzaussagen für gegebenes Niveau  $\alpha \in (0, 1)$ :

(i) Ist  $q_{F(p, n-p); 1-\alpha}$  das  $(1 - \alpha)$ -Quantil der  $F(p, n - p)$ -Verteilung, so ist

$$C := \{\beta \in \mathbb{R}^p \mid |X(\beta - \hat{\beta})|^2 \leq p \hat{\sigma}^2 q_{F(p, n-p); 1-\alpha}\}$$

ein Konfidenzellipsoid zum Konfidenzniveau  $1 - \alpha$  für  $\beta$ .

(ii) Ist  $q_{t(n-p); 1-\alpha/2}$  das  $(1 - \frac{\alpha}{2})$ -Quantil der  $t(n - p)$ -Verteilung, so ist

$$I := \left[ \hat{\rho} - \hat{\sigma} \sqrt{\langle (X^\top X)^{-1} v, v \rangle} q_{t(n-p); 1-\alpha/2}, \hat{\rho} + \hat{\sigma} \sqrt{\langle (X^\top X)^{-1} v, v \rangle} q_{t(n-p); 1-\alpha/2} \right]$$

ein Konfidenzintervall zum Konfidenzniveau  $1 - \alpha$  für  $\rho = \langle v, \beta \rangle$ .

Eine allgemeinere Klasse statistischer Fragestellungen im linearen Modell sind lineare (bzw. affine) Testprobleme. Sie bieten eine Vielzahl von Anwendungsmöglichkeiten und sind aufgrund ihrer hohen Relevanz standardmäßig in Statistik-Software, wie zum Beispiel R (siehe `linear.hypothesis` im `car`-Paket), implementiert.

**Definition 2.45** Im gewöhnlichen linearen Modell ist ein (zweiseitiges) **lineares Testproblem** gegeben durch

$$H_0: K\beta = c, \sigma > 0 \quad \text{gegen} \quad H_1: K\beta \neq c, \sigma > 0$$

für eine (deterministische) Matrix  $K \in \mathbb{R}^{r \times p}$  mit vollem Rang  $\text{rank}(K) = r \leq p$  und einem Vektor  $c \in \mathbb{R}^r$ .  $K$  wird **Kontrastmatrix** genannt.

Für eine bessere Übersichtlichkeit verzichten wir im Folgenden auf die Bedingung  $\sigma > 0$  in der Angabe von Hypothese und Alternative. Unter der Hypothese  $H_0$  werden also insgesamt  $r \leq p$  linear unabhängige Bedingungen an die Parameter des linearen Modells gestellt. Neben den Testproblemen die mit dem F-Test aus Satz 2.40 und dem t-Test aus Satz 2.42 behandelt wurden, umfassen lineare Testprobleme weitere wichtige Beispiele.

*Beispiel 2.46 (Lineare Testprobleme)*

1. Ein Test auf Gleichheit zweier Regressionskoeffizienten ist für  $j, l \in \{1, \dots, p\}$ ,  $j \neq l$ , gegeben durch

$$H_0: \beta_j = \beta_l \quad \text{gegen} \quad H_1: \beta_j \neq \beta_l.$$

Dies wird durch die Kontrastmatrix  $K = (a_{1,i}) \in \mathbb{R}^{1 \times p}$  mit  $a_{1,i} = \mathbb{1}_{\{i=j\}} - \mathbb{1}_{\{i=l\}}$  und  $c = 0$  modelliert.

2. Der Globaltest

$$H_0: \forall j \in \{1, \dots, d\} : \beta_j = 0 \quad \text{gegen} \quad H_1: \exists j \in \{1, \dots, d\} : \beta_j \neq 0$$

wird mit der Kontrastmatrix  $K = E_p$  und  $c = (0, \dots, 0)^\top$  beschrieben.

3. Der Test eines Subvektors  $\beta^* = (\beta_1^*, \dots, \beta_r^*)^\top$  mit  $r \leq p$ , das heisst

$$H_0: \forall j \in \{1, \dots, r\} : \beta_j = \beta_j^* \quad \text{gegen} \quad H_1: \exists j \in \{1, \dots, r\} : \beta_j \neq \beta_j^*,$$

führt auf die Kontrastmatrix  $K = (\mathbb{1}_{i=j})_{i,j} \in \mathbb{R}^{r \times p}$  und  $c = \beta^*$ .

Die Grundidee für das Testen linearer Hypothesen ist es, die Residuen des auf  $H_0: K\beta = c$  eingeschränkten Kleinst-Quadrat-Schätzers mit denen des uneingeschränkten Kleinst-Quadrat-Schätzers  $\hat{\beta}$  zu vergleichen.

**Definition 2.47** Der auf die Hypothese  $H_0: K\beta = c$  **eingeschränkte Kleinst-Quadrat-Schätzer**  $\hat{\beta}_{H_0}$  ist gegeben durch

$$\hat{\beta}_{H_0} := \arg \min_{\beta \in \mathbb{R}^p : K\beta = c} |Y - X\beta|^2. \quad (2.8)$$

Wir bezeichnen die **Summe der quadrierten Residuen** von  $\hat{\beta}$  und  $\hat{\beta}_{H_0}$  mit

$$RSS = |Y - X\hat{\beta}|^2 \quad \text{bzw.} \quad RSS_{H_0} := |Y - X\hat{\beta}_{H_0}|^2.$$

Per Definitionem ist  $RSS_{H_0} \geq RSS$ . Ist die Abweichung zu groß, spricht dies gegen die Hypothese  $H_0$ . In der Tat führt auch die Methode des Likelihood-Quotiententests zu diesem Ansatz.

**Lemma 2.48** *Im gewöhnlichen linearen Modell unter Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  besitzt jeder nichtrandomisierte Likelihood-Quotiententest für*

$$H_0: K\beta = c \quad \text{gegen} \quad H_1: K\beta \neq c$$

mit Kontrastmatrix  $K \in \mathbb{R}^{r \times p}$  und  $c \in \mathbb{R}^r$  die Form

$$\varphi_\alpha = \mathbb{1}(F > \tilde{c}_\alpha) \quad \text{mit} \quad F := \frac{n-p}{r} \frac{RSS_{H_0} - RSS}{RSS}, \quad \tilde{c}_\alpha \geq 0.$$

**Beweis** Beachtet man, dass  $\hat{\beta}_{H_0}$  der Maximum-Likelihood-Schätzer unter der Nullhypothese ist, so ergibt sich genau wie im Beweis von Satz 2.40 für die Likelihood-Quotientenstatistik

$$\begin{aligned} T &= \frac{\sup_{\beta \in \mathbb{R}^p, \sigma > 0} L(\beta, \sigma)}{\sup_{\beta \in \mathbb{R}^p \text{ mit } K\beta = c, \sigma > 0} L(\beta, \sigma)} = \left( \frac{|Y - X\hat{\beta}_{H_0}|^2}{|Y - X\hat{\beta}|^2} \right)^{n/2} \\ &= \left( 1 + \frac{RSS_{H_0} - RSS}{RSS} \right)^{n/2}. \end{aligned}$$

Durch monotone Transformation lässt sich daher ein nichtrandomisierter Likelihood-Quotiententest als  $\varphi_\alpha = \mathbb{1}(F > \tilde{c}_\alpha)$  mit  $\tilde{c}_\alpha \geq 0$  schreiben.  $\square$

Um den kritischen Wert  $\tilde{c}_\alpha$  von  $\varphi_\alpha$  zu bestimmen, müssen wir die Teststatistik  $F$ , auch *Fisher-Statistik* genannt, genauer analysieren. Wir erhalten folgende umfassende Aussage:

**Satz 2.49 (F-Test für lineare Hypothesen)** *Im gewöhnlichen linearen Modell unter Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  ist die lineare Hypothese*

$$H_0: K\beta = c \quad \text{gegen} \quad H_1: K\beta \neq c$$

mit Kontrastmatrix  $K \in \mathbb{R}^{r \times p}$  und  $c \in \mathbb{R}^r$  zu testen.

(i) Für den Schätzer  $\widehat{\beta}_{H_0}$  aus (2.8) gilt

$$\widehat{\beta}_{H_0} = \widehat{\beta} - (X^\top X)^{-1} K^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - c). \quad (2.9)$$

(ii) Es gilt

$$RSS_{H_0} - RSS = |X(\widehat{\beta} - \widehat{\beta}_{H_0})|^2 = \langle (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - c), K\widehat{\beta} - c \rangle$$

und unter  $H_0$  ist  $(RSS_{H_0} - RSS)/\sigma^2 \chi^2(r)$ -verteilt.

(iii) Der **F-Test**

$$\varphi_\alpha = \mathbb{1}(F > q_{F(r, n-p), 1-\alpha}) \text{ mit } F := \frac{n-p}{r} \frac{RSS_{H_0} - RSS}{RSS}$$

ist ein Likelihood-Quotiententest zum Niveau  $\alpha \in (0, 1)$ .

**Beweis** (i) Wir bezeichnen die rechte Seite von (2.9) mit  $\widetilde{\beta}_{H_0}$  und müssen zeigen, dass  $\widetilde{\beta}_{H_0}$  die eindeutige Lösung der Optimierung in (2.8) ist. Zunächst weisen wir die Nebenbedingung nach. In der Tat gilt:

$$K\widetilde{\beta}_{H_0} = K\widehat{\beta} - K(X^\top X)^{-1} K^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - c) = c,$$

Aus Lemma 2.14 wissen wir, dass  $Y - X\widehat{\beta} = (E_n - \Pi_X)Y \perp \text{Im}(X)$ , sodass für  $\gamma \in \mathbb{R}^p$  nach dem Satz von Pythagoras

$$|Y - X\gamma|^2 = |Y - X\widehat{\beta} + X(\widehat{\beta} - \gamma)|^2 = |Y - X\widehat{\beta}|^2 + |X(\widehat{\beta} - \gamma)|^2$$

gilt. Außerdem ist

$$|X(\widehat{\beta} - \gamma)|^2 = |X(\widehat{\beta} - \widetilde{\beta}_{H_0})|^2 + |X(\widetilde{\beta}_{H_0} - \gamma)|^2 + 2\langle X(\widehat{\beta} - \widetilde{\beta}_{H_0}), X(\widetilde{\beta}_{H_0} - \gamma) \rangle.$$

Die Wahl von  $\widetilde{\beta}_{H_0}$  impliziert jedoch für  $\gamma$  mit  $K\gamma = c$

$$\begin{aligned} & \langle X(\widehat{\beta} - \widetilde{\beta}_{H_0}), X(\widetilde{\beta}_{H_0} - \gamma) \rangle \\ &= ((X^\top X)^{-1} K^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - c))^\top X^\top X(\widetilde{\beta}_{H_0} - \gamma) \\ &= (K\widehat{\beta} - c)^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widetilde{\beta}_{H_0} - K\gamma) = 0, \end{aligned}$$

wobei die letzte Gleichheit aus  $K\widetilde{\beta}_{H_0} = K\gamma = c$  folgt. Insgesamt erhalten wir

$$|Y - X\gamma|^2 = |Y - X\widehat{\beta}|^2 + |X(\widehat{\beta} - \widetilde{\beta}_{H_0})|^2 + |X(\widetilde{\beta}_{H_0} - \gamma)|^2, \quad (2.10)$$

was offensichtlich für  $\gamma = \tilde{\beta}_{H_0}$  minimal ist, sodass  $\hat{\beta}_{H_0} = \tilde{\beta}_{H_0}$  gilt.

(ii) Aus (2.10) mit  $\gamma = \hat{\beta}_{H_0}$  folgt durch Einsetzen von  $\hat{\beta}_{H_0}$

$$\begin{aligned} RSS_{H_0} - RSS &= |Y - X\hat{\beta}_{H_0}|^2 - |Y - X\hat{\beta}|^2 = |X(\hat{\beta} - \hat{\beta}_{H_0})|^2 \\ &= (\hat{\beta} - \hat{\beta}_{H_0})^\top X^\top X (\hat{\beta} - \hat{\beta}_{H_0}) \\ &= (K\hat{\beta} - c)^\top (K(X^\top X)^{-1}K^\top)^{-1} (K\hat{\beta} - c). \end{aligned}$$

Nach dem Satz von Gauß-Markov ist  $\hat{\beta}$  unter  $H_0$  ein erwartungstreuer Schätzer von  $\beta$  mit  $\text{Cov}_0(\hat{\beta}) = \sigma^2(X^\top X)^{-1}$ . Es ergibt  $H_0$

$$\mathbb{E}_0[K\hat{\beta} - c] = K\mathbb{E}_0[\hat{\beta}] - c = K\beta - c = 0,$$

und

$$\text{Var}_0(K\hat{\beta} - c) = K \text{Cov}_0(\hat{\beta}) K^\top = \sigma^2 K(X^\top X)^{-1} K^\top.$$

Aus der Normalverteilung von  $\hat{\beta}$  folgt daher  $(RSS_{H_0} - RSS)/\sigma^2 \sim \chi^2(r)$ .

(iii) Da  $RSS_{H_0} - RSS$  eine messbare Funktion von  $\hat{\beta}$  und somit auch von  $X\hat{\beta} = \Pi_X Y$  ist, ist  $RSS = |(E_n - \Pi_X)Y|^2$  unabhängig von  $RSS_{H_0} - RSS$ .  $F \sim F(r, n - p)$  folgt daher aus der Charakterisierung der  $F(r, n - p)$ -Verteilung aus Lemma 2.37. Wir schließen mit Lemma 2.48.  $\square$

*Bemerkung 2.50*  $W := rF$  heißt auch *Wald-Statistik*. Unter Verwendung von Satz 2.49(ii) und Lemma 2.18 können wir die Fisher-Statistik auch als

$$F = \frac{\frac{1}{r}|X\hat{\beta} - X\hat{\beta}_{H_0}|^2}{\hat{\sigma}^2}$$

schreiben.

*Beispiel 2.51 (t-Test als Spezialfall des F-Tests)* Wir betrachten den Spezialfall einer eindimensionalen Nebenbedingung. Mit  $r = 1$ ,  $K = v^\top$ ,  $\rho = \langle v, \beta \rangle$ ,  $c = \rho_0$  testen wir also  $H_0: \rho = \rho_0$  gegen  $H_1: \rho \neq \rho_0$ . Satz 2.49(ii) zeigt mit  $\hat{\rho} = \langle v, \hat{\beta} \rangle$

$$F = \frac{|X\hat{\beta} - X\hat{\beta}_{H_0}|^2}{\hat{\sigma}^2} = \frac{(\hat{\rho} - \rho_0)^2}{\hat{\sigma}^2 \langle (X^\top X)^{-1}v, v \rangle} \sim F(1, n - p) \text{ unter } H_0.$$

Damit ist der  $F$ -Test äquivalent zum zweiseitigen t-Test mit der Teststatistik

$$T = \frac{\hat{\rho} - \rho_0}{\hat{\sigma} \sqrt{\langle (X^\top X)^{-1}v, v \rangle}} \sim t(n - p) \text{ unter } H_0.$$

Wir erhalten in diesem Fall also genau Satz 2.42.

Wir wollen nun die entwickelte Theorie auf reale Daten anwenden.

*Beispiel 2.52 (Klimadaten)* Wir betrachten die mittleren Julitemperaturen zwischen 1719 und 2020, die an einer Wetterstation in Berlin-Dahlem gemessen wurden. Bis

auf wenige Ausnahmen haben wir Messungen aus jedem Jahr. Insgesamt liegen  $n = 291$  Beobachtungen vor, die über das Data Climate Center des Deutschen Wetterdienstes<sup>1</sup> frei verfügbar sind. Eine polynomiale Regression in der Zeit  $t$  (in Jahrhunderten beginnend bei 1719, siehe Abbildung 2.6) mit Polynomgraden  $d = 1, \dots, 4$  liefert (mit gerundeten Werten)

$$\begin{aligned} p_1(t) &= 18,72 - 0,002 t, \\ p_2(t) &= 18,63 + 0,16 t - 0,05 t^2, \\ p_3(t) &= 17,43 + 4,49 t - 3,52 t^2 + 0,75 t^3, \\ p_4(t) &= 17,39 + 4,71 t - 3,84 t^2 + 0,92 t^3 - 0,03 t^4. \end{aligned}$$

Wir sehen hier sogar einen leichten negativen Anstieg in  $p_1$ , was den allgemeinen Erkenntnissen zur Klimaentwicklung zu widersprechen scheint. Andererseits ist der Faktor  $-0,002$  nur sehr klein. Es fällt zudem auf, dass auch der jeweils letzte Grad in  $p_2$  und  $p_4$  einen nur sehr geringen Einfluss auf die Regressionskurve hat. Welches Polynom ist nun am besten geeignet, um die Daten zu beschreiben?

Um diese Frage beantworten zu können, verwenden wir die vorangegangene Testtheorie. Zunächst ist es plausibel, dass die zufälligen Schwankungen zwischen den Jahren unabhängig voneinander sind und als näherungsweise normalverteilt angenommen werden können (QQ-Plot). Um statistisch verwertbare Aussagen zu treffen, setzen wir noch das Niveau  $\alpha = 0,05$  fest. Der Parametervektor ist  $\beta = (\beta_0, \dots, \beta_d)^\top$ .

*Frage 1:* Ist der negative Trend von  $p_1$  signifikant, wenn wir das lineare Modell mit  $d = 1$  annehmen?  $H_0: \beta_1 \geq 0$  gegen  $H_1: \beta_1 < 0$ . Die zugehörige t-Statistik  $T = -\frac{\hat{\beta}_1}{\hat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v}} \approx 0,01$  mit  $v = (0, 1)^\top \in \mathbb{R}^2$  liegt deutlich unter dem kritischen Wert  $q_{t(n-2), 1-\alpha} \approx 1,65$  (einseitiger T-Test), sodass die Hypothese nicht verworfen werden kann. Es gibt also keinen signifikant negativen Trend.

*Frage 2:* Liegt den Beobachtungen (im Modell mit  $d = 4$ ) ein linearer Zusammenhang zugrunde?  $H_0: \beta_2 = \beta_3 = \beta_4 = 0$ . Mittels Bemerkung 2.50 (oder direkt über die quadrierten Residuen) berechnen wir die Fisher-Statistik

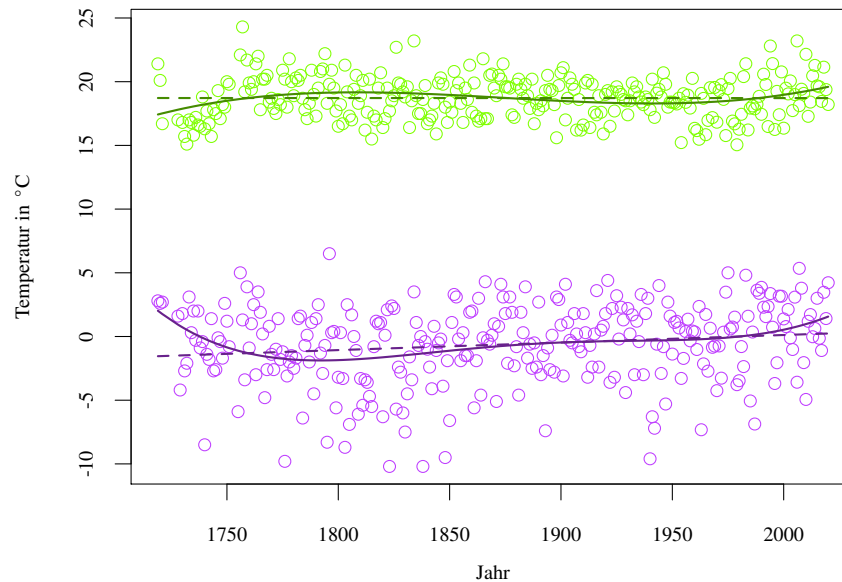
$$F = \frac{\sum_{i=1}^n (p_4(t_i) - p_1(t_i))^2}{3\hat{\sigma}^2} \approx 5,25 > 2,63 \approx q_{F(3, n-5), 1-\alpha}.$$

Folglich kann die Hypothese abgelehnt werden, und wir schlussfolgern, dass eine Regressionsgerade unzureichend ist.

*Frage 3:* Benötigen wir ein Polynom vierten Grades?  $H_0: \beta_4 = 0$ . Die zugehörige t-Statistik hat den Wert  $-0,11$ , dessen Absolutbetrag kleiner als das Quantil  $q_{t(n-5); 0,975} \approx 1,97$  ist (zweiseitiger t-Test). Diese Nullhypothese kann also akzeptiert werden.

*Frage 4:* Benötigen wir ein Polynom dritten Grades?  $H_0: \beta_3 = 0$  (im Modell mit  $d = 3$ ). Die zugehörige t-Statistik hat den Wert  $3,96$ , dessen Absolutbetrag größer

<sup>1</sup> Siehe <https://cdc.dwd.de/portal/>



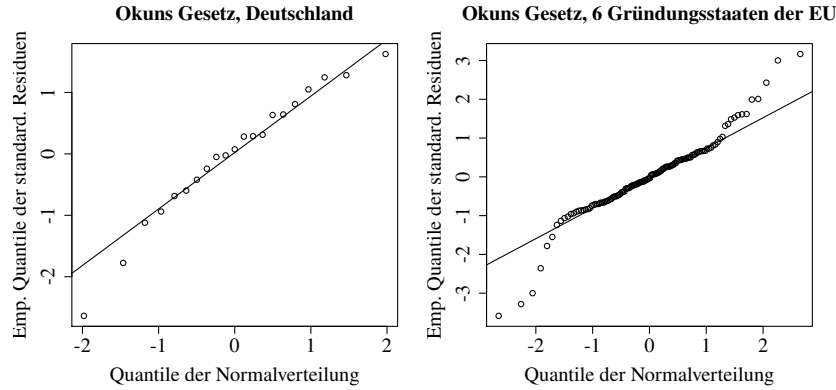
**Abb. 2.6** Polynomielle Regression über die mittleren Julitemperaturen (grün) und Januartemperaturen (violett) in Berlin-Dahlem von 1719 bis 2020: Gestrichelte Linien zeigen die jeweiligen Regressionsgraden. Durchgezogene Linien zeigen das Regressionpolynom dritten Grades für den Juli und vierten Grades für den Januar. Datenbasis: Deutscher Wetterdienst

als das Quantil  $q_{t(n-4);0,975} \approx 1,97$  ist. Die Hypothese kann also abgelehnt werden, und der kubische Anteil im Regressionspolynom ist signifikant.

Es sei darauf hingewiesen, dass wir die jeweiligen Testprobleme einzeln betrachten. Wenn wir alle Fragen simultan zu einem Niveau  $\alpha$  beantworten wollen, so liegt ein multiples Testproblem vor, und die kritischen Werte müssen korrigiert werden (sogenannte *Multiplizitätskorrektur*), beispielsweise indem man bei  $m$  simultanen Tests jeweils das Niveau  $\alpha/m$  verwendet (*Bonferroni-Korrektur*).

$p_3$  zeigt einen deutlichen Anstieg der Temperaturen in der zweiten Hälfte des 20. Jahrhunderts. Da wir hier nur eine einzelne Zeitreihe betrachtet haben, kann daraus aber kein allgemeiner Zusammenhang geschlossen werden. Das muss Ergebnis einer Kooperation mit Klimatologen sein. Eine analoge Analyse der Januar-Mitteltemperaturen zeigt übrigens einen signifikanten Koeffizienten vierten Grades (Aufgabe 2.17) und ebenfalls einen deutlichen Anstieg in den letzten 100 Jahren.

Zum Abschluss dieses Kapitels sei nochmal betont, dass vor der Anwendung der Inferenztheorie aus diesem Kapitel in Praxis zunächst geprüft werden muss, ob ein gewöhnliches lineares Modell unter Normalverteilungsannahme tatsächlich geeignet ist, um die beobachteten Daten  $(x_i, Y_i)_{i=1, \dots, n}$  zu beschreiben.



**Abb. 2.7** QQ-Plots der standardisierten Residuen im einfachen linearen Modell für das jährliche Wachstum des BIP in Abhängigkeit von der jährlichen Veränderung der Arbeitslosenquote für Deutschland (links) und die 6 Gründungsstaaten der EU (rechts) in den Jahren 1992 bis 2012

Um dies zu prüfen, wird unter anderem der aus Lemma 2.18 bekannte Residuenvektor

$$R = Y - X\hat{\beta} = (E_n - \Pi_X)\varepsilon$$

herangezogen, wobei  $\hat{\beta}$  der Kleinste-Quadrate-Schätzer und  $\Pi_X = X(X^\top X)^{-1}X^\top$  die Projektionsmatrix auf  $\text{Im } X$  sind. Unter den Modellannahmen sind die Residuen  $R_i$  zentriert und normalverteilt, aber nicht mehr unabhängig voneinander, und sie besitzen verschiedene Varianzen  $(1 - h_i)\sigma^2$  mit  $h_i = (X(X^\top X)^{-1}X^\top)_{ii}$ . Ein Plot der Residuen gegen die vorhergesagten Werte  $X\hat{\beta} = \Pi_X Y$  gibt Aufschluss, ob ein systematischer Modellfehler vorliegt. Hier sollte keine Abhängigkeit erkennbar sein, was aus Aufgabe 2.8 folgt.

Um die Normalverteilungsannahme zu überprüfen, bietet sich ein QQ-Plot der standardisierten Residuen

$$T_i = \frac{R_i}{\sqrt{(1 - h_i)\hat{\sigma}^2}}$$

an. Man beachte, dass wir die unbekannte Varianz  $\sigma^2$  durch ihren Schätzer aus Lemma 2.18 ersetzt haben, sodass  $T_i$  nur noch approximativ normalverteilt ist.

**Beispiel 2.53 (QQ-Plots für Okuns Gesetz)** Wir greifen den empirischen Zusammenhang zwischen der Änderung der Arbeitslosenquote  $x_i$  und Wachstum des Bruttoinlandsproduktes  $Y_i$  aus Beispiel 2.2 auf. Abbildung 2.7 zeigt die aus dem einfachen linearen Modell  $Y_i = ax_i + b + \varepsilon_i$  für die Jahre  $i = 1992, \dots, 2012$  resultierenden QQ-Plots, wobei Deutschland einzeln und zusammen mit 5 weiteren EU-Staaten betrachtet wird. In beiden Fällen sehen wir, dass die Normalverteilungsannahme im Zentrum gut zutrifft, jedoch sehr große und sehr kleine Quantile zum Teil zu deutlichen Abweichungen führen.

Statistische Tests zum Überprüfen der Verteilungsannahme beruhen beispielsweise auf dem  $\chi^2$ -Test oder dem Kolmogorov-Smirnov-Test, auf die wir hier jedoch nicht eingehen werden. Die interessierte Leserin sei auf Lehmann and Romano (2005) verwiesen. Die gesamte Modellverifikation für das lineare Modell wird von Fahrmeir et al. (2009) ausführlich diskutiert.

## 2.3 Varianzanalyse

In der Varianzanalyse werden Varianzen verschiedener Gruppen benutzt, um mögliche Unterschiede zwischen diesen Gruppen nachzuweisen. Die Varianzanalyse beruht auf dem Testen von Spezialfällen linearer Hypothesen im linearen Modell. Wir werden also unsere bisherigen Resultate anwenden, aber gewisse Strukturen zusätzlich ausnutzen.

*Beispiel 2.54 (Düngemittel – a)* Um den Einfluss von  $p \in \mathbb{N}$  verschiedenen Düngemitteln auf den Ernteertrag zu vergleichen, wird jedes Düngemittel  $i \in \{1, \dots, p\}$  auf  $n_i$  verschiedenen Agrarflächen ausgebracht. Der durch Witterungseinflüsse etc. zufällige Ernteertrag kann mittels  $Y_{ij} = \mu_i + \varepsilon_{ij}$  für  $j = 1, \dots, n_i$  und  $i = 1, \dots, p$  modelliert werden, wobei  $\mu_i$  der mittlere Ernteertrag von Düngemittel  $i$  ist und  $\varepsilon_{ij}$  unabhängige, zentrierte Störgrößen sind. Die zu untersuchende Frage ist, ob die einzelnen Düngemittel einen unterschiedlichen Einfluss auf den mittleren Ernteertrag haben.

**Definition 2.55** Das Modell der **einfaktoriellen Varianzanalyse** (englisch: *(one-way) analysis of variance*, kurz: ANOVA1) ist gegeben durch Beobachtungen

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad i = 1, \dots, p, j = 1, \dots, n_i,$$

mit i.i.d.-verteilten Störgrößen  $\varepsilon_{ij} \sim N(0, \sigma^2)$ . Wir bezeichnen den ersten Index als den **Faktor** und den Wert  $i = 1, \dots, p$  als die **Faktorstufe**. Folglich geben  $(n_i)_{i=1, \dots, p}$  die Anzahl der unabhängigen Versuchswiederholungen pro Faktor an und  $n := \sum_{i=1}^p n_i$  ist der Gesamtstichprobenumfang. Gilt  $n_1 = \dots = n_p$ , so sprechen wir von **balanciertem Design**.

*Beispiel 2.56 (Düngemittel – b)* Der Faktor sind die Düngemittel, und weil wir zwei Düngemittel testen, haben wir zwei Faktorstufen  $i \in \{1, 2\}$ . Der Gesamtstichprobenumfang beträgt  $n_1 + n_2 = 2 + 3 = 5$ , und es liegt wegen  $n_1 \neq n_2$  kein balanciertes Design vor.

*Bemerkung 2.57 (ANOVA1)* Das ANOVA1-Modell ist ein Spezialfall des gewöhnlichen linearen Modells der Form

$$\mathbb{R}^n \ni Y := \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ \vdots \\ Y_{k1} \\ \vdots \\ Y_{kn_p} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & & & \vdots \\ 1 & 0 & \cdots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \cdots & 1 \\ \vdots & & & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}}_{=: X \in \mathbb{R}^{n \times p}} \cdot \underbrace{\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}}_{=: \mu \in \mathbb{R}^p} + \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \vdots \\ \varepsilon_{p1} \\ \vdots \\ \varepsilon_{pn_p} \end{pmatrix}.$$

Hierbei gilt  $\text{Im}(X) = p$ .

Die klassische Fragestellung der Varianzanalyse lautet: „Existieren Unterschiede in den faktorstufenspezifischen Mittelwerten  $\mu_i$ ?“ oder anders formuliert „Hat der Faktor einen Einfluss auf die Response-Variable oder nicht?“. Dies führt auf das *Testproblem*

$$H_0: \mu_1 = \cdots = \mu_p \quad \text{gegen} \quad H_1: \exists i, l \in \{1, \dots, p\} : \mu_i \neq \mu_l. \quad (2.11)$$

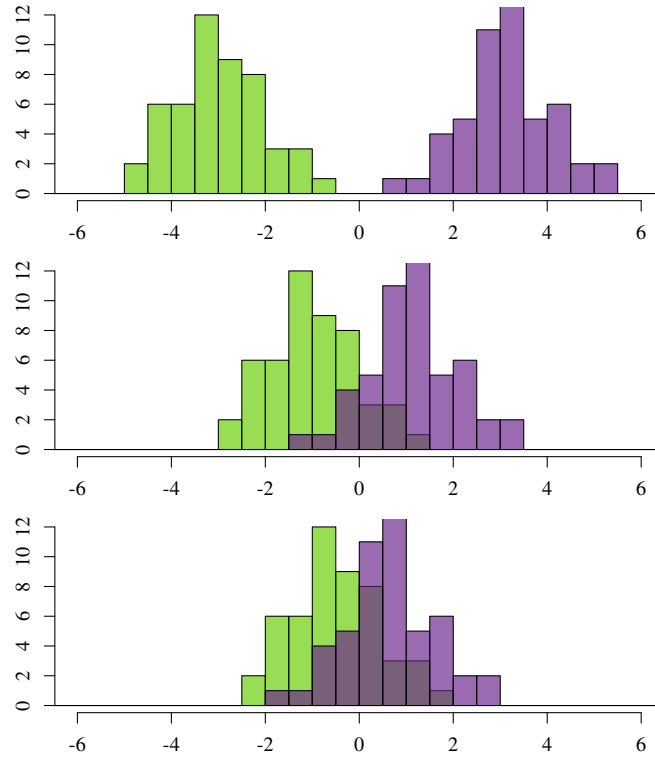
Eine erste Idee, um diese Hypothese zu überprüfen, wäre die Mittelwerte jeder Faktorstufe zu berechnen und diese zu vergleichen. Eine große Abweichung würde gegen die Hypothese sprechen. Dieser Vergleich muss jedoch die Streuung der Beobachtungen um den jeweiligen Mittelwert berücksichtigen, da andernfalls nicht klar ist, was eine große Abweichung ist (man vergleiche beispielsweise  $N(0; 1)$  mit  $N(0, 1; 1)$  und  $N(0; 0, 01)$  mit  $N(0, 1; 0, 01)$ ). Wir sollten also *die Nullhypothese ablehnen, wenn die Streuung zwischen den Gruppen größer ist als die Streuung innerhalb der Gruppen*. Dieses Vorgehen motiviert die Bezeichnung Varianzanalyse.

*Beispiel 2.58 (Düngemittel – c)* Wir bleiben bei zwei Faktorstufen (zwei Düngemittel), erhöhen aber zur besseren Illustration die Anzahl der Felder pro Düngemittel  $n_1 = n_2 = 50$ . Abbildung 2.8 zeigt die Histogramme von drei verschiedenen Ernteszenarien, bei denen sich die Mittelwerte annähern, während die Varianzen gleich bleiben. Man sieht deutlich, dass im ersten Fall beide Faktorstufen leicht voneinander zu unterscheiden sind, während dies im dritten Fall kaum noch möglich ist.

**Lemma 2.59 (Streuungszerlegung)** *Im Modell der einfaktoriellen Varianzanalyse definieren wir das  $i$ -te Gruppenmittel,  $i = 1, \dots, p$ , bzw. das Gesamtmittel als*

$$\bar{Y}_{i\bullet} := \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad \text{bzw.} \quad \bar{Y}_{\bullet\bullet} := \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^{n_i} Y_{ij}$$

*sowie die Streuungsmaße*



**Abb. 2.8** Teilweise überlappende Histogramme von drei verschiedenen, normalverteilten Beobachtungen von zwei Gruppen (grün und violett). Von oben nach unten rücken die Erwartungswerte der Gruppen näher aneinander (erst liegen sie bei  $\pm 3$ , dann bei  $\pm 1$  und zuletzt bei  $\pm 0,5$ ) mit gleichbleibenden Varianzen ( $\sigma^2 = 1$ ).

$$SSB := \sum_{i=1}^p n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 \quad \text{und} \quad SSW := \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2.$$

Dann gilt

$$SST := \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet})^2 = SSB + SSW.$$

**Beweis** Es gilt

$$\begin{aligned} SST &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{\bullet\bullet})^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{i\bullet} + \bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 \\ &= \sum_i \sum_j ((Y_{ij} - \bar{Y}_{i\bullet})^2 + 2(Y_{ij} - \bar{Y}_{i\bullet})(\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) + (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2), \end{aligned}$$

wobei

$$\begin{aligned}
\sum_i \sum_j (Y_{ij} - \bar{Y}_{i\bullet})(\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) &= \sum_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) \sum_j (Y_{ij} - \bar{Y}_{i\bullet}) \\
&= \sum_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})(n_i \bar{Y}_{i\bullet} - n_i \bar{Y}_{i\bullet}) = 0.
\end{aligned}$$

Damit ist die Darstellung von  $SST$  gezeigt.  $\square$

**Bemerkung 2.60** Nach Normierung mit  $1/n$  handelt es sich bei den Größen  $SSB$ ,  $SSW$  bzw.  $SST$  um die gewichteten empirischen Varianzen der Mittelwerte der Gruppen (englisch: *sum of squares between groups*), der Summe der empirischen Varianz innerhalb der Gruppen (englisch: *sum of squares within groups*) bzw. der empirischen Varianz der gesamten Stichprobe (englisch: *total sum of squares*).

**Satz 2.61** Im Modell der einfaktoriellen Varianzanalyse mit  $n > p \geq 2$  gilt:

(i) Der Kleinste-Quadrate-Schätzer von  $\mu = (\mu_1, \dots, \mu_p)^\top$  ist gegeben durch

$$\hat{\mu} = (\bar{Y}_{1\bullet}, \dots, \bar{Y}_{p\bullet})^\top.$$

- (ii)  $SSW/\sigma^2 \sim \chi^2(n-p)$  und unter  $H_0: \mu_1 = \dots = \mu_p$  gilt  $SSB/\sigma^2 \sim \chi^2(p-1)$ .  
(iii)  $SSW$  und  $SSB$  sind unabhängig und somit  $F := \frac{n-p}{p-1} \frac{SSB}{SSW} \stackrel{H_0}{\sim} F(p-1, n-p)$ .

**Beweis** (i) Nachrechnen zeigt

$$\hat{\mu} = (X^\top X)^{-1} X^\top Y = \begin{pmatrix} 1/n_1 & & 0 \\ & \ddots & \\ 0 & & 1/n_p \end{pmatrix} \begin{pmatrix} \sum_{j=1}^{n_1} Y_{1j} \\ \vdots \\ \sum_{j=1}^{n_p} Y_{pj} \end{pmatrix} = \begin{pmatrix} \bar{Y}_{1\bullet} \\ \vdots \\ \bar{Y}_{p\bullet} \end{pmatrix}.$$

(ii) und (iii) Wegen  $SSW = |Y - X\hat{\mu}|^2 = |R|^2$  für die Residuen  $R$  aus Lemma 2.18 folgt  $SSW/\sigma^2 \sim \chi^2(n-p)$  und die Unabhängigkeit von  $SSW$  und  $\hat{\mu}$  aus Lemma 2.39. Nach dem vorangegangenen Satz gilt weiterhin  $SSB = SST - SSW$ . Somit folgt die Behauptung aus Satz 2.49, falls  $SST = |Y - X\hat{\mu}_{H_0}|^2$ , wobei  $\hat{\mu}_{H_0}$  gegeben ist durch

$$|Y - X\hat{\mu}_{H_0}|^2 = \min_{\mu \in \mathbb{R}} \underbrace{\left| Y - X \begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix} \right|^2}_{\in \mathbb{R}^p} = \min_{\mu \in \mathbb{R}} \underbrace{\left| Y - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \mu \right|^2}_{=: X_0 \in \mathbb{R}^{n \times 1}}.$$

Dieses Minimierungsproblem wird durch  $\hat{\mu}_{H_0} = (X_0^\top X_0)^{-1} X_0^\top Y = n^{-1} \sum_{i,j} Y_{ij} = \bar{Y}_{\bullet\bullet}$  gelöst. Damit folgt die Behauptung.  $\square$

Folglich können wir die Hypothese aus (2.11) mit dem Likelihood-Quotiententest aus Satz 2.49, das heißt einem F-Test, überprüfen.

	FG	Quadratsummen	Quadratmittel	F-Statistik
zwischen	$p - 1$	$SSB = \sum_{i=1}^p n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2$	$SSB/(p - 1)$	$\frac{n - p}{p - 1} \frac{SSB}{SSW}$
innerhalb	$n - p$	$SSW = \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2$	$SSW/(n - p)$	
total	$n - 1$	$SST = \sum_{i=1}^p \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet})^2$	$SST/(n - 1)$	

**Tabelle 2.1** ANOVA1-Tafel zur Darstellung von Freiheitsgraden (FG), Quadratsummen, Quadratmittel und der resultierenden F-Statistik in der einfaktoriellen Varianzanalyse

	FG	Quadratsummen	Quadratmittel	F-Statistik
zwischen	1	$SSB = 199,314$	$SSB/1 = 199,314$	$F = 8131,686$
innerhalb	98	$SSW = 2,402$	$SSW/98 = 0,025$	
total	99	$SST = 201,716$	$SST/99 = 2,037$	
	FG	Quadratsummen	Quadratmittel	F-Statistik
zwischen	1	$SSB = 7,471$	$SSB/1 = 7,471$	$F = 86,700$
innerhalb	98	$SSW = 8,445$	$SSW/98 = 0,086$	
total	99	$SST = 15,916$	$SST/99 = 0,161$	
	FG	Quadratsummen	Quadratmittel	F-Statistik
zwischen	1	$SSB = 12,171$	$SSB/1 = 12,171$	$F = 10,281$
innerhalb	98	$SSW = 116,014$	$SSW/98 = 1,184$	
total	99	$SST = 128,185$	$SST/99 = 1,295$	

**Tabelle 2.2** ANOVA1-Tafeln für den Ernteertrag unter dem Einsatz von zwei verschiedenen Düngemitteln in drei verschiedenen Szenarien

**Methode 2.62 (Einfaktorielle Varianzanalyse)** Im Modell der einfaktoriellen Varianzanalyse testen wir

$$H_0: \mu_1 = \dots = \mu_p \quad \text{versus} \quad H_1: \exists i, l \in \{1, \dots, p\} : \mu_i \neq \mu_l$$

zum Niveau  $\alpha \in (0, 1)$  durch den F-Test

$$\varphi_\alpha = \mathbb{1}(F > q_{F(p-1, n-p); 1-\alpha}) \quad \text{mit} \quad F := \frac{n-p}{p-1} \frac{SSB}{SSW}.$$

Es ist übersichtlich, die einzelnen Zwischenergebnisse der ANOVA1 in eine Tabelle zu schreiben (siehe Tabelle 2.1 und das folgende Beispiel).

**Beispiel 2.63 (Düngemittel – d)** Wie zuvor betrachten wir zwei Düngemittel, die jeweils auf 50 Feldern eingesetzt werden und wenden den F-Test auf die drei Beobachtungssätze aus Abbildung 2.8 an. Für die drei Szenarien erhalten wir ANOVA1-Tafeln, siehe Tabelle 2.2.

Es gilt  $q_{F(1,98);0,95} = 3,938$ ,  $q_{F(1,98);0,99} = 6,901$  und  $q_{F(1,98);0,999} = 11,510$ . Bis auf den Test zum Niveau  $\alpha = 0.001$  der dritten Beobachtung wird die Nullhypothese also immer abgelehnt.

**Bemerkung 2.64 (Effektdarstellung)** Das einfaktorielle Varianzanalysemodell lässt sich zu

$$Y_{ij} = \mu_0 + \alpha_i + \varepsilon_{ij}, \quad i = 1, \dots, p, j = 1, \dots, n_i,$$

umformen, wobei  $\mu_0 := \frac{1}{n} \sum_{i=1}^p n_i \mu_i = \mathbb{E}[\bar{Y}_{\bullet\bullet}]$  das Gesamtmittel ist und  $\alpha_i := \mu_i - \mu_0$ ,  $i = 1, \dots, p$ , den **Effekt der Faktorstufe** beschreibt. Diese Form heißt **Effektdarstellung**, und sie verlangt die Nebenbedingung

$$0 = \sum_{i=1}^p n_i \alpha_i \quad \text{oder äquivalent} \quad \alpha_p = -\frac{1}{n_p} \sum_{i=1}^{p-1} n_i \alpha_i,$$

damit die Designmatrix weiter vollen Rang hat. Der Parametervektor ist also gegeben durch  $(\mu_0, \alpha_1, \dots, \alpha_{p-1})^\top$ . Die F-Statistik, um die Globalhypothese

$$H_0: \alpha_1 = \dots = \alpha_{p-1} = 0$$

zu überprüfen, ist identisch zur Statistik aus Satz 2.61. Per Konstruktion lässt sich somit anhand der Schätzungen des Parametervektors ablesen, wie stark der Effekt, zum Beispiel von Düngemitteln, im Vergleich zum Gesamtdurchschnitt ist (negativ wie auch positiv).

Im Fall  $p = 2$  führt die Varianzanalyse auf den Zweistichproben-t-Test.

**Korollar 2.65 (Zweistichproben-t-Test)** Im Modell der einfaktoriellen Varianzanalyse mit  $p = 2$  und dem Testproblem  $H_0: \mu_1 = \mu_2, \sigma > 0$  gegen  $H_1: \mu_1 \neq \mu_2, \sigma > 0$  ist

$$\varphi_\alpha = \mathbb{1}(|T| > q_{t(n-2), 1-\alpha/2}) \quad \text{mit} \quad T := \frac{\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet}}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})SSW/(n-2)}}$$

mit dem  $(1 - \alpha/2)$ -Quantil  $q_{t(n-2), 1-\alpha/2}$  der  $t(n-2)$ -Verteilung der sogenannte **Zweistichproben-t-Test** der Hypothese  $H_0$  zum Niveau  $\alpha \in (0, 1)$ .

**Beweis** Wegen  $p = 2$  und Satz 2.61 gilt unter  $H_0$

$$\frac{n-p}{p-1} \cdot \frac{SSB}{SSW} = \frac{SSB}{SSW/(n-2)} \sim F(1, n-2).$$

Zur Erinnerung gilt  $T_n^2 = F_{1,n}$ . Wegen  $n\bar{Y}_{\bullet\bullet} = n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet}$  gilt

$$\begin{aligned} SSB &= n_1(\bar{Y}_{1\bullet} - \bar{Y}_{\bullet\bullet})^2 + n_2(\bar{Y}_{2\bullet} - \bar{Y}_{\bullet\bullet})^2 \\ &= n_1\bar{Y}_{1\bullet}^2 + n_2\bar{Y}_{2\bullet}^2 + n\bar{Y}_{\bullet\bullet}^2 - 2(n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet})\bar{Y}_{\bullet\bullet} \\ &= n_1\bar{Y}_{1\bullet}^2 + n_2\bar{Y}_{2\bullet}^2 - \frac{1}{n}(n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet})^2 = \frac{n_1n_2}{n}(\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet})^2. \end{aligned}$$

Folglich gilt unter  $H_0$

$$T := \frac{\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet}}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})SSW/(n-2)}} \sim t(n-2).$$

**Bemerkung 2.66 (Zweistichproben-t-Test)** Alternativ kann die Teststatistik des Zweistichproben-t-Tests auch als

$$T = \frac{1}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \frac{\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet}}{\sqrt{\frac{n_1-1}{n-2}\hat{\sigma}_1^2 + \frac{n_2-1}{n-2}\hat{\sigma}_2^2}} \quad \text{mit} \quad \hat{\sigma}_i^2 = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$$

geschrieben werden. Dabei wird die Varianz  $\sigma^2$  durch das gewichtete Mittel der gruppenweisen empirischen Varianzen  $\hat{\sigma}_i^2$  geschätzt.

Sind die Varianzen in den Grundgesamtheiten (zum Beispiel die Felder mit den verschiedenen Düngemitteln) ungleich, dann kann obiges Resultat nicht angewendet werden. Eine Modifikation des Zweistichproben-t-Test führt auf den Welch-Test, der den verschiedenen Varianzen Rechnung trägt, siehe Lehmann and Romano (2005).

Nach dem Studium der einfaktoriellen Varianzanalyse liegt es nahe, den Einfluss von mehreren Faktoren zu berücksichtigen. Wir beschränken uns hier auf zwei Faktoren. Eine Varianzanalyse mit mehr als zwei Faktoren führt zu ähnlichen Verteilungsaussagen und resultierenden Tests.

**Definition 2.67** Das Modell der **zweifaktoriellen Varianzanalyse** (kurz: ANOVA2) mit balanciertem Design ist gegeben durch Beobachtungen

$$\begin{aligned} Y_{ijk} &= \mu_{ij} + \varepsilon_{ijk}, \quad i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K \\ &= \mu_0 + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk} \end{aligned}$$

mit  $I, J, K \geq 2$ , i.i.d.-verteilten Störgrößen  $\varepsilon_{ijk} \sim N(0, \sigma^2)$  und Nebenbedingungen (der Effektdarstellung)

$$\sum_{i=1}^I \alpha_i = \sum_{j=1}^J \beta_j = \sum_{i=1}^I \gamma_{ij} = \sum_{j=1}^J \gamma_{ij} = 0.$$

Wir haben also zwei Faktoren mit Faktorstufen  $i = 1, \dots, I$  und  $j = 1, \dots, J$ .  $(\alpha_i)$  und  $(\beta_j)$  heißen **Haupteffekte** des ersten beziehungsweise zweiten Faktors.  $(\gamma_{ij})$  heißen **Interaktions-** oder **Wechselwirkungseffekte**.

Das ANOVA2-Modell ist also ein lineares Modell mit zwei kategoriellen Kovariablen. Die Gesamtanzahl an Beobachtungen ist gegeben durch  $n = I \cdot J \cdot K$  und die Parameterdimension ist  $p = I \cdot J$ . Die typische Testprobleme sind

$$H_0: \quad \forall i: \alpha_i = 0 \quad \text{gegen} \quad H_1: \exists i \in \{1, \dots, I\}: \alpha_i \neq 0, \quad (2.12)$$

$$H_0: \quad \forall j: \beta_j = 0 \quad \text{gegen} \quad H_1: \exists j \in \{1, \dots, J\}: \beta_j \neq 0, \quad (2.13)$$

$$H_0: \forall i, j: \gamma_{ij} = 0 \quad \text{gegen} \quad H_1: \exists i \in \{1, \dots, I\}, j \in \{1, \dots, J\}: \gamma_{ij} \neq 0. \quad (2.14)$$

**Satz 2.68** Im Modell der zweifaktoriellen Varianzanalyse mit balanciertem Design gilt:

(i) Die Kleinsten-Quadrate-Schätzer für  $\mu_0, \alpha_i, \beta_j$  und  $\gamma_{ij}$ ,  $i = 1, \dots, I-1, j = 1, \dots, J-1$ , sind gegeben durch

$$\begin{aligned}\hat{\mu}_0 &= \bar{Y}_{\bullet\bullet\bullet}, & \hat{\alpha}_i &= \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet\bullet\bullet}, & \hat{\beta}_j &= \bar{Y}_{\bullet j\bullet} - \bar{Y}_{\bullet\bullet\bullet}, \\ \hat{\gamma}_{ij} &= (\bar{Y}_{ij\bullet} - \bar{Y}_{\bullet\bullet\bullet}) - \hat{\alpha}_i - \hat{\beta}_j = \bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j\bullet} + \bar{Y}_{\bullet\bullet\bullet},\end{aligned}$$

wobei wir wieder in den mit  $\bullet$  gekennzeichneten Koordinaten Mittelwerte bilden.

(ii) Definieren wir

$$\begin{aligned}SSW &:= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij\bullet})^2, \\ SSB_1 &:= JK \sum_{i=1}^I (\bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet\bullet\bullet})^2, & SSB_2 &:= IK \sum_{j=1}^J (\bar{Y}_{\bullet j\bullet} - \bar{Y}_{\bullet\bullet\bullet})^2, \\ SSB_{12} &:= K \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j\bullet} + \bar{Y}_{\bullet\bullet\bullet})^2,\end{aligned}$$

dann können die Hypothesen (2.12), (2.13) bzw. (2.14) mit den F-Statistiken

$$\begin{aligned}\frac{IJ(K-1)}{I-1} \frac{SSB_1}{SSW} &\sim F(I-1, IJ(K-1)), \\ \frac{IJ(K-1)}{J-1} \frac{SSB_2}{SSW} &\sim F(J-1, IJ(K-1)) \quad \text{bzw.} \\ \frac{IJ(K-1)}{(I-1)(J-1)} \frac{SSB_{12}}{SSW} &\sim F((I-1)(J-1), IJ(K-1))\end{aligned}$$

getestet werden, wobei die Statistiken jeweils unter der Nullhypothese F-verteilt sind.

Den Beweis überlassen wir den Leserinnen und Lesern als Übung (Aufgabe 2.18).

**Beispiel 2.69 (ANOVA2)** Eine Bäuerin interessiert sich neben dem Effekt ihrer Düngemittel auch für den Einfluss der genutzten Samenarten. Sie verwendet  $I = 3$  verschiedene Samenarten und  $J = 5$  verschiedene Düngemittel. Das ergibt  $3 \cdot 5 = 15$  verschiedene Feldergruppen, auf denen je genau eine Samenart und ein Düngemittel ausgebracht wird. Jede Samen-Düngemittel-Kombination wird auf  $K = 2$  Feldern eingesetzt. Die Erträge der Bäuerin sind in Tabelle 2.3 zusammengefasst.

Für  $\alpha = 0,05$  erhalten wir die ANOVA-Tafel in Tabelle 2.4. Die Hypothesen, dass die Haupteffekte null seien, werden also verworfen, aber die Hypothese, dass die Interaktionseffekte null sind, wird bestätigt.

	DM1	DM2	DM3	DM4	DM5
SA1	111;116	100;106	99;113	108;110	105;108
SA2	115;118	103;105	105;107	113;118	110;113
SA3	99;103	91;93	103;105	104;107	99;104

**Tabelle 2.3** Messung der Ernteerträge von je zwei Feldern pro Samenart-Düngemittel- Kombination. DM steht für Düngemittel und SA für Samenart.

	FG	Quadratsummen	Quadratmittel	F-Statistik	F-Quantil
Zwischen	2	$SSB_1 = 512,867$	256,433	20,298	3,682
	4	$SSB_2 = 449,467$	112,367	8,894	3,056
	8	$SSB_{12} = 143,133$	17,892	1,416	2,641
Innerhalb	15	$SSW = 189,500$	12,633		
Total	29	$SST = 1294,967$	44,654		

**Tabelle 2.4** ANOVA2-Tafel, wobei in der ersten Spalte von rechts das  $(1 - \alpha)$ -Quantil der F-Verteilung mit entsprechenden Freiheitsgraden steht

## 2.4 Aufgaben

### 2.1 Bestimmen Sie im einfachen linearen Modell

$$Y_i = ax_i + b + \varepsilon_i \quad \text{für } i = 1, \dots, n, n \in \mathbb{N}$$

die Maximum-Likelihood-Schätzer für die unbekannten Parameter  $a, b \in \mathbb{R}$  aufgrund der Beobachtungen  $Y_1, \dots, Y_n$ , wenn

- $(\varepsilon_1, \dots, \varepsilon_n)$  ein Vektor unabhängiger (zentrierter) Laplace-verteilter Zufallsvariablen mit Skalierungsparameter  $\beta > 0$  ist, das heißt die Verteilung von  $\varepsilon_i$  hat die Lebesgue-Dichte  $f_\beta(x) = \frac{2}{\beta} e^{-|x|/\beta}$ ,  $x \in \mathbb{R}$ ,
- $(\varepsilon_1, \dots, \varepsilon_n)$  ein Vektor unabhängiger Exp( $\lambda$ )-verteilter Zufallsvariablen mit Parameter  $\lambda > 0$  ist, das heißt die Verteilung von  $\varepsilon_i$  hat die Lebesgue-Dichte  $g_\lambda(x) = \lambda e^{-\lambda x} \mathbb{1}_{[0, \infty)}(x)$ ,  $x \in \mathbb{R}$ .

### 2.2 Untersuchen Sie den Zusammenhang zwischen Bruttoinlandsprodukt und Happiness-Score aus Beispiel 1.3. Bestimmen Sie anhand der Daten aus dem Jahr 2019 die Regressionsgrade. Recherchieren Sie das Bruttoinlandsprodukt von Deutschland im Jahr 2020. Welchen Happiness-Score würden Sie vorhersagen?

### 2.3 Zehn erkrankte Patienten nehmen dasselbe Medikament, jedoch in verschiedenen Dosen. Folgende Tabelle zeigt die Anzahl an Tagen bis zur Genesung:

- Verwenden Sie ein einfaches lineares Modell, um festzustellen, ob eine höhere Dosis zu einer schnelleren Genesung führt.

- (b) Die ersten fünf Patienten in der Tabelle waren Frauen, während die hinteren fünf Patienten Männer waren. Schätzen Sie für jede der beiden Gruppen eine eigene Regressionsgrade. ändert dies die Schlussfolgerung aus (a)?

Dosis	45	55	70	60	75	80	100	90	110	125
Genesungszeit	5	6	8	8	9	3	4	6	5	7

- 2.4 Im linearen Modell  $Y = X\beta + \varepsilon$  mit dem Mittelwert der Beobachtung  $\bar{Y}$  und den geschätzten Zielgrößen  $\hat{Y} := X\hat{\beta}$  heißt

$$R^2 := \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

*Bestimmtheitsmaß.* Zeigen Sie, dass  $R^2$  in  $[0, 1]$  liegt. Weisen Sie nach, dass  $R^2 = 0$  einen linearen Zusammenhang ausschließt und dass  $R^2 = 1$  einen perfekt linearen Zusammenhang impliziert.

- 2.5 Leiten Sie aus dem Satz von Gauß-Markov ab, dass die optimale Varianz im gewöhnlichen linearen Modell

- (a) von  $\hat{\beta}$  gleich  $\sigma^2(X^T X)^{-1}$  ist,  
 (b) von  $\langle v, \hat{\beta} \rangle$  gleich  $\sigma^2 |X(X^T X)^{-1} v|^2$  ist.

- 2.6 Bestimmen Sie im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$ ,  $\sigma > 0$ , den Maximum-Likelihood-Schätzer  $\hat{\sigma}_{ML}^2$  von  $\sigma^2$ . Vergleichen Sie Bias und Varianz von  $\hat{\sigma}_{ML}^2$  und dem Schätzer  $\hat{\sigma}^2 = |Y - X\hat{\beta}|^2 / (n - k)$  aus Lemma 2.18.

- 2.7 Beweisen Sie im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$ ,  $\sigma > 0$ , dass für den Kleinste-Quadrate-Schätzer  $\hat{\beta}$  gilt:

$$\frac{1}{n} \mathbb{E}[|X\beta - X\hat{\beta}|^2] = \sigma^2 \frac{p}{n}.$$

- 2.8 Im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$ ,  $\sigma > 0$ , und Kleinste-Quadrate-Schätzer  $\hat{\beta} = X^\dagger Y$  sind  $\hat{Y} := (\hat{Y}_1, \dots, \hat{Y}_n)^\top := X\hat{\beta}$  die geschätzten Erwartungswerte und  $R = (R_1, \dots, R_n)^\top = Y - \hat{Y}$  die Residuen. Beweisen Sie folgende geometrische Eigenschaften:

- (a)  $\hat{Y}$  ist orthogonal zu  $\hat{\varepsilon}$ , das heißt  $\langle \hat{Y}, R \rangle = 0$ .  
 (b) Die Spalten von  $X$  sind orthogonal zu den Residuen, d. h.  $X^\top R = 0$ .

Gilt zusätzlich, dass wir ein Regressionsmodell mit Absolutglied haben, das heißt in der ersten Spalte der Designmatrix  $X$  stehen nur Einsen, gilt weiterhin:

- (c) Die Residuen sind im Mittel gleich null, das heißt  $\sum_{i=1}^n R_i = 0$ .
- (d) Der arithmetische Mittelwert der  $\widehat{Y}_i$  ist gleich dem Mittelwert der Beobachtungen  $Y_i$  selbst, das heißt  $\sum_{i=1}^n \widehat{Y}_i = \sum_{i=1}^n Y_i$ .

2.9 Im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\beta \in \mathbb{R}^p$ ,  $\varepsilon \sim N(0, \sigma^2 E_n)$  und  $\sigma > 0$  bezeichne  $\widehat{\beta}$  den Kleinst-Quadrate-Schätzer und  $\widehat{\beta}_{\text{ridge}}$  den Ridge-Regressionsschätzer.

- (a) Finden Sie für den Ridge-Regressionsschätzer den optimalen Tuning-Parameter  $\lambda$  in Abhängigkeit vom unbekannten  $\beta \in \mathbb{R}^p$  und bestimmen Sie das resultierende minimale Risiko  $\mathbb{E}[|\widehat{\beta}_{\text{ridge}} - \beta|^2]$ .
- (b) Folgern Sie, dass es immer ein  $\lambda > 0$  gibt, sodass gilt

$$\mathbb{E}[|\widehat{\beta}_{\text{ridge}} - \beta|^2] < \mathbb{E}[|\widehat{\beta} - \beta|^2].$$

2.10 Untersuchen Sie das Verhalten des Ridge-Regressionsschätzers im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  mit  $\varepsilon \sim N(0, E_n)$  bei  $n = 50$  Beobachtungen und  $p = 30$  unbekannten Parametern, wenn davon 10 groß und 20 klein sind. Gehen Sie wie folgt vor:

- (a) Simulieren Sie die Einträge der Designmatrix  $X \in \mathbb{R}^{n \times p}$  als unabhängige  $N(0, 1)$ -verteilte Zufallsvariablen und erzeugen Sie den Vektor  $\beta \in \mathbb{R}^p$  aus 10  $U([\frac{1}{2}, 1])$ -verteilten und 20  $U([0, \frac{3}{10}])$ -verteilten Zufallsvariablen.
- (b) Erzeugen Sie in 200 Durchgängen jeweils den Fehlervektor

$$\varepsilon^{(i)} \in \mathbb{R}^{50}, i = 1, \dots, 200,$$

und berechnen Sie aus den resultierenden Beobachtungen den Ridge-Regressionsschätzer  $\widehat{\beta}_{\lambda}^{(i)}$  für  $\lambda \in \{\frac{k}{2} : k = 0, \dots, 50\}$ .

- (c) Bestimmen Sie für jedes  $\lambda$  den (empirischen) mittleren quadratischen Fehler  $R_{\lambda} := \frac{1}{200} \sum_{i=1}^{200} |\widehat{\beta}_{\lambda}^{(i)} - \beta|^2$ . Stellen Sie die Abbildung  $\lambda \mapsto R_{\lambda}$  graphisch dar und vergleichen Sie die Fehler der Ridge-Regressionsschätzer mit dem des gewöhnlichen Kleinst-Quadrate-Schätzers.

2.11 Beweisen Sie für die empirische Kovarianzmatrix  $\Sigma_n = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$  im Fall unabhängiger  $p$ -dimensionaler Zufallsvektoren  $X_i \sim N(0, \Sigma)$ :

- (a) Es gilt  $\mathbb{E}[\Sigma_n] = \Sigma$  (der Erwartungswert einer Matrix ist die Matrix der Erwartungswerte).
- (b) Mit der *Frobenius-Norm*  $\|M\|_2 = (\sum_{i,j=1}^p M_{i,j}^2)^{1/2}$  und Spur  $\text{tr}(M) = \sum_{i=1}^p M_{i,i}$  einer  $p \times p$ -Matrix  $M$  gilt

$$\mathbb{E}[\|\Sigma_n - \Sigma\|_2^2] = n^{-1}(\text{tr}(\Sigma)^2 + \|\Sigma\|_2^2), \quad \mathbb{E}[\|\Sigma^{-1/2}(\Sigma_n - \Sigma)\Sigma^{-1/2}\|_2^2] = n^{-1}(p^2 + p),$$

wobei in der zweiten Gleichung  $\Sigma$  als invertierbar angenommen wird.

- 2.12 Betrachten Sie das lineare Modell  $Y_i = X_i\beta + \varepsilon_i, i = 1, \dots, n$ , mit i.i.d.  $(X_i, Y_i)_{i=1, \dots, n} \subseteq [-R, R]^p \times \mathbb{R}$  für ein  $R > 0$  und  $\mathbb{E}[\varepsilon_i^2] = \sigma^2 > 0$ . Zudem seien  $X_i$  und  $\varepsilon_i$  für alle  $i$  unabhängig und  $\Sigma_X = \mathbb{E}[X_1 X_1^\top] \in \mathbb{R}^{p \times p}$  sei wohldefiniert und positiv definit. Beweisen Sie für den Kleinste-Quadrate-Schätzer  $\hat{\beta}$  die asymptotische Normalität:

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma^2 \Sigma_X^{-1}).$$

- 2.13 Betrachten Sie das gewöhnliche lineare Modell unter Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  mit unbekanntem  $\sigma > 0$  und wahren Parameter  $\beta \in \mathbb{R}^p$ . Für  $v, \beta_0 \in \mathbb{R}^p$  seien  $\rho = \langle v, \beta \rangle$  und  $\rho_0 = \langle v, \beta_0 \rangle$  gegeben. Zeigen Sie, dass der Likelihood-Quotiententest der Hypothese  $H_0: \rho = \rho_0, \sigma > 0$  gegen die Alternative  $H_1: \rho \neq \rho_0, \sigma > 0$  zum Niveau  $\alpha \in (0, 1)$  von der Form

$$\varphi_\alpha = \mathbb{1}(|T_{n-p}(Y)| > c_\alpha) \quad \text{mit} \quad T_{n-p}(Y) := \frac{\hat{\rho} - \rho_0}{\hat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v}}$$

und einem geeigneten kritischen Wert  $c_\alpha > 0$  ist.

- 2.14 Im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  unter der Normalverteilungsannahme  $(\varepsilon_1, \dots, \varepsilon_n)^\top \sim N(0, \sigma^2 E_n)$  mit unbekanntem  $\beta \in \mathbb{R}^p$  und  $\sigma > 0$  bestimmt der Kovariablenvektor  $x_{n+1} \in \mathbb{R}^p$  die zukünftige Beobachtung  $Y_{n+1} = \langle x_{n+1}, \beta \rangle + \varepsilon_{n+1}$  mit  $\varepsilon_{n+1} \sim N(0, \sigma^2)$  unabhängig von  $(\varepsilon_1, \dots, \varepsilon_n)$ . Sei  $\alpha \in (0, 1)$ .
- Konstruieren Sie ein  $(1 - \alpha)$ -Konfidenzintervall für den zu erwartenden Wert  $\langle x_{n+1}, \beta \rangle$ .
  - Konstruieren Sie ein  $(1 - \alpha)$ -Prognoseintervall für die zu beobachtende Realisierung von  $Y_{n+1}$ .

Geben Sie im Modell aus Aufgabe 2 beide Intervalle zum Niveau 0,95 für den Happiness-Score in Deutschland 2020 explizit an.

- 2.15 Zeigen Sie im gewöhnlichen linearen Modell  $Y = X\beta + \varepsilon$  unter der Normalverteilungsannahme  $(\varepsilon_1, \dots, \varepsilon_n)^\top \sim N(0, \sigma^2 E_n)$  mit unbekanntem  $\beta \in \mathbb{R}^k$  und  $\sigma > 0$ , dass

$$C := \left[ (n-k)\hat{\sigma}^2 / q_{\chi^2(n-k), 1-\alpha/2}, (n-k)\hat{\sigma}^2 / q_{\chi^2(n-k), \alpha/2} \right]$$

ein Konfidenzintervall für  $\sigma^2$  zum Konfidenzniveau  $1 - \alpha$  ist, wobei  $\hat{\sigma}^2 = |Y - X\hat{\beta}|^2 / (n-k)$  die Stichprobenvarianz und  $q_{\chi^2(n-k), \tau}$  für  $\tau \in (0, 1)$  das  $\tau$ -Quantil der  $\chi^2(n-k)$ -Verteilung sind. Nutzen Sie dieses Resultat, um einen  $\chi^2$ -Test für die Varianz  $\sigma^2$  zum Niveau  $\alpha \in (0, 1)$  zu konstruieren.

- 2.16 Im gewöhnlichen linearen Modell unter Normalverteilungsannahme  $\varepsilon \sim N(0, \sigma^2 E_n)$  mit  $\sigma^2 > 0$  soll die lineare Hypothese

$$H_0 : \beta_j = \beta_l \quad \text{gegen} \quad H_1 : \beta_j \neq \beta_l.$$

getestet werden. Zeigen Sie, dass die zugehörige Fisher-Statistik von der Form

$$F = \frac{(\hat{\beta}_j - \hat{\beta}_l)^2}{\widehat{\text{Var}}(\hat{\beta}_j - \hat{\beta}_l)}$$

und unter  $H_0$   $F(1, n - p)$ -verteilt ist, wobei  $\widehat{\text{Var}}(\hat{\beta}_j - \hat{\beta}_l)$  ein geeigneter Schätzer der Varianz  $\text{Var}(\hat{\beta}_j - \hat{\beta}_l)$  ist. Warum ist dieser F-Test äquivalent zum (zweiseitigen) t-Test mit der Teststatistik

$$T = \frac{\hat{\beta}_j - \hat{\beta}_l}{(\widehat{\text{Var}}(\hat{\beta}_j - \hat{\beta}_l))^{1/2}} \stackrel{H_0}{\sim} t(n - p)?$$

2.17 Führen Sie eine lineare Regressionsanalyse analog zu Beispiel 2.52 für die mittleren Januartemperaturen in Berlin-Dahlem zwischen 1719 und 2020 basierend auf den Messungen des Deutschen Wetterdienstes aus.

2.18 Beweisen Sie im Modell der zweifaktoriellen Varianzanalyse den Satz 2.68.



## Kapitel 3

# Modellwahl

Am Ende des vorangegangenen Kapitels haben wir eine mögliche Modellmisspezifikation betrachtet, die Daten werden also durch eine Verteilung erzeugt, die gar nicht in dem Modell enthalten ist, das unserem statistischen Verfahren zugrunde liegt. Dies wirft die Frage auf, was überhaupt ein geeignetes Modell ist. In diesem Kapitel werden wir verschiedene Ansätze kennenlernen, die eine versierte Wahl des Modells ermöglichen, nämlich die Informationskriterien AIC und BIC, die Idee der Kreuzvalidierung und die Lasso-Methode. Um diese Methoden zu analysieren, werden wir Orakelungleichungen herleiten.

### 3.1 Informationskriterien

Im Allgemeinen wird ein statistisches Modell die zugrunde liegenden Zusammenhänge nicht exakt darstellen. Unabhängig davon, ob wir die möglicherweise sehr komplexen Abhängigkeiten kennen oder nicht, ist eine exakte Beschreibung oft gar nicht das Ziel. Stattdessen möchten wir ein möglichst einfaches Modell, das sowohl die gegebenen Beobachtungen als auch neue Daten aus demselben Experiment möglichst gut beschreibt. Bei der Wahl eines Modells müssen daher zwei entgegengesetzt wirkende Probleme ausbalanciert werden:

1. Ist das Modell zu stark vereinfacht, werden Haupteffekte nicht mehr ausreichend erklärt. Dieses sogenannte *underfitting* führt zu einer systematischen Verzerrung/einem Bias.
2. Ist das Modell sehr komplex mit vielen Parametern, dann können die gegebenen Beobachtungen sehr genau dargestellt werden, bishin zur Interpolation der Daten. Viele Modellparameter führen jedoch zu einem großen statistischen Fehler und starken Schwankungen, sodass man für neue Beobachtungen keine vernünftige Modellvorhersage mehr erwarten kann. Man spricht von *overfitting*.

Folgendes Beispiel verdeutlicht diese beiden Effekte:

*Beispiel 3.1* Für i.i.d. zentrierte Fehler  $(\varepsilon_i)_{i=1,\dots,n}$  beobachten wir

$$Y_i = f(x_i) + \varepsilon_i, \quad 1 \leq i \leq n, \quad (3.1)$$

mit einer unbekannten Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  und  $x_i \in \mathbb{R}^d$ ,  $i = 1, \dots, n$ . Wir betrachten nun die empirische Seminorm

$$\|f\|_n := \sqrt{\langle f, f \rangle_n}, \quad \text{wobei} \quad \langle f, g \rangle_n := \frac{1}{n} \sum_{i=1}^n f(x_i)g(x_i),$$

auf dem Funktionenraum  $\mathcal{F} := \{f: \mathbb{R}^d \rightarrow \mathbb{R}\}$ . Für eine Orthonormalbasis  $(e_j)_{j=1,\dots,n}$  dieses Funktionenraumes bezüglich  $\langle \cdot, \cdot \rangle_n$  können wir die Regressionsfunktion durch

$$f(x_i) = \sum_{j=1}^n \langle f, e_j \rangle_n e_j(x_i), \quad i = 1, \dots, n,$$

darstellen. Statt ein allgemeines  $f \in \mathcal{F}$  zu betrachten, können wir hoffen, dass die unbekannte Regressionsfunktion  $f$  für ein  $p \in \{1, \dots, n\}$  bereits durch ein Element aus dem linearen Unterraum  $\text{span}\{e_1, \dots, e_p\}$  gut approximiert werden kann. Diese Annahme führt auf das vereinfachte lineare Modell

$$Y_i = \sum_{j=1}^p \beta_j e_j(x_i) + \varepsilon_i, \quad 1 \leq i \leq n, \quad p \in \mathbb{N},$$

sodass die Modellparameter  $(\beta_1, \dots, \beta_p) \in \mathbb{R}^p$  durch ein  $\widehat{\beta}^{(p)} = (\widehat{\beta}_1, \dots, \widehat{\beta}_p)$  geschätzt werden müssen. Da die Daten jedoch gemäß (3.1) erzeugt werden, liegt hier eine *Modellmisspezifikation* vor (siehe auch Kapitel ??). Dennoch kann ein resultierender Schätzer  $\widehat{f}_{n,p} = \sum_{j=1}^p \widehat{\beta}_j e_j$  eine gute Näherung von  $f$  sein.

Wir haben somit verschiedene Modelle mit Parameterdimensionen  $p = 1, \dots, n$  zur Auswahl. Wie groß sollen wir  $p$  wählen? Ein großes  $p$  bedeutet, dass der (quadrierte) Bias  $\|f - \sum_{j=1}^p \beta_j e_j\|_n^2 = \sum_{j=p+1}^n \langle f, e_j \rangle^2$  klein ist. Gleichzeitig führt die große Parameterdimension zu großen stochastischen Fehlern, denn wir müssen  $p$  Koeffizienten schätzen, siehe (2.3) und Aufgabe 2.7. Wir sollten also  $p$  so wählen, dass Bias und stochastischer Fehler annähernd gleich sind.

In diesem Kapitel wollen wir Methoden entwickeln, die automatisch und anhand der gegebenen Daten ein geeignetes Modell wählen. Hierzu untersuchen wir Daten aus einem Wahrscheinlichkeitsraum  $(\mathcal{X}, \mathcal{F}, \mathbb{P})$ , wobei wir das zugrunde liegende Wahrscheinlichkeitsmaß  $\mathbb{P}$  nicht kennen. Wie in Beispiel 3.1 illustriert, modellieren wir die Beobachtungen durch  $P \in \mathbb{N}$  verschiedene parametrische Modelle

$$(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\theta)_{\theta \in \Theta_p}),$$

wobei die Dimension  $p \in \{1, \dots, P\}$  des Parameterraums  $\Theta_p$  wächst und  $P \in \mathbb{N}$  die maximal betrachtete Parameterdimension ist. Innerhalb des  $p$ -ten Modells liefert

uns die Maximum-Likelihood-Methode einen Schätzer  $\hat{\vartheta}_p$ . Unser Ziel ist, ein  $\hat{p}$  zu finden, sodass  $\mathbb{P}_{\hat{\vartheta}_p}$  die wahre Verteilung  $\mathbb{P}$  gut approximiert. Man beachte dabei, dass  $\mathbb{P}$  in keiner der Familien  $(\mathbb{P}_{\vartheta})_{\vartheta \in \Theta_p}$  liegen muss.

Um zu entscheiden, ob ein Kandidat  $\mathbb{P}_{\hat{\vartheta}_p}$  die unbekannte Verteilung  $\mathbb{P}$  gut approximiert, benötigen wir ein Abstandsmaß zwischen den Wahrscheinlichkeitsverteilungen.

**Definition 3.2** Für Wahrscheinlichkeitsmaße  $\mathbb{P}$  und  $\mathbb{Q}$  auf  $(X, \mathcal{F})$  heißt

$$\text{KL}(\mathbb{P}|\mathbb{Q}) = \begin{cases} \int_X \log\left(\frac{d\mathbb{P}}{d\mathbb{Q}}(x)\right) \mathbb{P}(dx), & \text{falls } \mathbb{P} \ll \mathbb{Q}, \\ \infty, & \text{sonst,} \end{cases}$$

**Kullback-Leibler-Divergenz** von  $\mathbb{P}$  bezüglich  $\mathbb{Q}$ .

Besitzen  $\mathbb{P}$  und  $\mathbb{Q}$   $\mu$ -Dichten  $p$  und  $q$ , so berechnet man  $\text{KL}(\mathbb{P}|\mathbb{Q}) = \int_{\{q>0\}} \log(p/q) p d\mu$  im Fall  $\mathbb{P} \ll \mathbb{Q}$ . Die Kullback-Leibler-Divergenz ist keine Metrik auf dem Raum der Wahrscheinlichkeitsmaße (es gilt weder Symmetrie noch Dreiecksungleichung, vergleiche Aufgabe 4.1), jedoch besitzt sie nützliche Eigenschaften, um verschiedene Wahrscheinlichkeitsmaße zu vergleichen.

**Lemma 3.3 (Eigenschaften der Kullback-Leibler-Divergenz)** Für Wahrscheinlichkeitsmaße  $\mathbb{P}$  und  $\mathbb{Q}$  auf  $(X, \mathcal{F})$  gilt:

- (i)  $\text{KL}(\mathbb{P}|\mathbb{Q}) \geq 0$  und  $\text{KL}(\mathbb{P}|\mathbb{Q}) = 0$  gilt genau dann, wenn  $\mathbb{P} = \mathbb{Q}$ .
- (ii)  $\text{KL}(\mathbb{P}^{\otimes n}|\mathbb{Q}^{\otimes n}) = n \text{KL}(\mathbb{P}|\mathbb{Q})$  für alle  $n \in \mathbb{N}$ .
- (iii) Bildet  $(\mathbb{P}_{\eta})_{\eta \in \Xi}$  eine natürliche  $p$ -parametrische Exponentialfamilie der Form

$$\frac{d\mathbb{P}_{\eta}}{d\mu}(x) = c(x) \exp(\langle \eta, T(x) \rangle - \zeta(\eta))$$

und liegt  $\eta_0$  im Inneren von  $\Xi$ , so ist

$$\text{KL}(\mathbb{P}_{\eta_0}|\mathbb{P}_{\eta}) = \zeta(\eta) - \zeta(\eta_0) - \langle \nabla \zeta(\eta_0), \eta - \eta_0 \rangle, \quad \eta \in \Xi.$$

- (iv) Für unabhängige und gemäß  $\mathbb{P}$  identisch verteilte Beobachtungen  $(X_i)_{i \in \mathbb{N}}$  mit Werten in  $(X, \mathcal{F})$  und Loglikelihood-Funktion  $l(\vartheta) := \log \frac{d\mathbb{P}_{\vartheta}^{X_1}}{d\mu}$  in einem statistischen Modell bezüglich einem dominierenden Maß  $\mu$  auf  $(X, \mathcal{F})$  gilt

$$-\frac{1}{n} \sum_{i=1}^n l(\vartheta, X_i) \xrightarrow{\mathbb{P}\text{-f.s.}} -\mathbb{E}_{\mathbb{P}}[l(\vartheta, X_1)] = \text{KL}(\mathbb{P}|\mathbb{P}_{\vartheta}^{X_1}) - \text{KL}(\mathbb{P}|\mu) \quad \text{für } n \rightarrow \infty,$$

sofern die Kullback-Leibler-Divergenzen (für allgemeine Maße  $\mu$  analog definiert) endlich sind.

**Beweis**

- (i) Ohne Einschränkung können wir  $\mathbb{P} \ll \mathbb{Q}$  annehmen, da andernfalls  $\text{KL}(\mathbb{P}|\mathbb{Q}) = \infty > 0$  stets erfüllt ist. Dann gilt

$$\text{KL}(\mathbb{P}|\mathbb{Q}) = \int_X \underbrace{\log\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) \frac{d\mathbb{P}}{d\mathbb{Q}}}_{=: f\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right)} d\mathbb{Q}$$

für  $f(x) = x \log(x)$  mit  $f(0) = 0$ . Wegen  $f''(x) = \frac{1}{x} > 0$  ist  $f$  strikt konvex. Die Jensen-Ungleichung liefert somit

$$\text{KL}(\mathbb{P}|\mathbb{Q}) = \int_X f\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{Q} \geq f\left(\int_X \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{Q}\right) = f(1) = 0.$$

Hierbei gilt genau dann Gleichheit, wenn  $d\mathbb{P}/d\mathbb{Q}$   $\mathbb{Q}$ -f.s. konstant ist.

- (ii) Wir können wieder  $\mathbb{P} \ll \mathbb{Q}$  annehmen. Da die Radon-Nikodym-Dichte von Produktmaßen gleich dem Produkt der Randdichten ist, folgt

$$\begin{aligned} \text{KL}(\mathbb{P}^{\otimes n}|\mathbb{Q}^{\otimes n}) &= \int_{X^n} \log\left(\frac{d\mathbb{P}^{\otimes n}}{d\mathbb{Q}^{\otimes n}}\right) d\mathbb{P}^{\otimes n} \\ &= \sum_{i=1}^n \int_X \log\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{P} = n \text{KL}(\mathbb{P}|\mathbb{Q}). \end{aligned}$$

- (iii) Einsetzen der  $\mu$ -Dichten liefert

$$\begin{aligned} \text{KL}(\mathbb{P}_{\eta_0}|\mathbb{P}_{\eta}) &= \mathbb{E}_{\eta_0} \left[ \log \left( \frac{d\mathbb{P}_{\eta_0}/d\mu}{d\mathbb{P}_{\eta}/d\mu} \right) \right] \\ &= \mathbb{E}_{\eta_0} \left[ \log \left( \frac{d\mathbb{P}_{\eta_0}}{d\mu} \right) - \log \left( \frac{d\mathbb{P}_{\eta}}{d\mu} \right) \right] \\ &= \mathbb{E}_{\eta_0} [\langle \eta_0, T \rangle - \zeta(\eta_0) - (\langle \eta, T \rangle - \zeta(\eta))] \\ &= \zeta(\eta) - \zeta(\eta_0) - \langle \eta - \eta_0, \mathbb{E}_{\eta_0}[T] \rangle. \end{aligned}$$

Schließlich gilt  $\nabla \zeta(\eta_0) = \mathbb{E}_{\eta_0}[T]$  nach Satz ??(iii).

- (iv) Unter den Annahmen  $\text{KL}(\mathbb{P}|\mathbb{P}_{\vartheta}) < \infty$  und  $\text{KL}(\mathbb{P}|\mu) < \infty$  gilt

$$\begin{aligned} \text{KL}(\mathbb{P}|\mathbb{P}_{\vartheta}) - \text{KL}(\mathbb{P}|\mu) &= \int_X \left( \log \left( \frac{d\mathbb{P}}{d\mathbb{P}_{\vartheta}} \right) - \log \left( \frac{d\mathbb{P}}{d\mu} \right) \right) d\mathbb{P} \\ &= - \int_X \log \left( \frac{d\mathbb{P}_{\vartheta}}{d\mu} \right) d\mathbb{P}. \end{aligned}$$

Damit ist  $l(\vartheta) = \log\left(\frac{d\mathbb{P}_{\vartheta}}{d\mu}\right) \in \mathcal{L}^1(\mathbb{P})$ , und das starke Gesetz der großen Zahlen liefert für  $\mathbb{P}$ -verteilte, unabhängige  $(X_i)$

$$-\frac{1}{n} \sum_{i=1}^n l(\vartheta, X_i) \xrightarrow{\mathbb{P}\text{-f.s.}} -\mathbb{E}_{\mathbb{P}}[l(\vartheta, X_1)] = \text{KL}(\mathbb{P}|\mathbb{P}_{\vartheta}) - \text{KL}(\mathbb{P}|\mu)$$

für  $n \rightarrow \infty$ . □

Die Kullback-Leibler-Divergenz von Normalverteilungen ergibt sich als Spezialfall von (iii):

**Korollar 3.4** Für  $\mu_0, \mu \in \mathbb{R}^p$  und eine symmetrische, positiv-definite Matrix  $\Sigma \in \mathbb{R}^{p \times p}$  gilt

$$\text{KL}(\mathcal{N}(\mu_0, \Sigma) | \mathcal{N}(\mu, \Sigma)) = \frac{1}{2} |\Sigma^{-1/2}(\mu - \mu_0)|^2 = \frac{1}{2} \langle \Sigma^{-1}(\mu - \mu_0), \mu - \mu_0 \rangle.$$

**Beweis** Wir verwenden (iii) aus dem vorangegangenen Lemma, wobei  $\zeta(\mu) = \frac{1}{2} \langle \Sigma^{-1} \mu, \mu \rangle$ . Daraus folgt

$$\begin{aligned} \text{KL}(\mathcal{N}(\mu_0, \Sigma) | \mathcal{N}(\mu, \Sigma)) &= \frac{1}{2} \left( \langle \Sigma^{-1} \mu, \mu \rangle - \langle \Sigma^{-1} \mu_0, \mu_0 \rangle - 2 \langle \Sigma^{-1}(\mu - \mu_0), \mu_0 \rangle \right) \\ &= \frac{1}{2} \left( \langle \Sigma^{-1} \mu, \mu \rangle + \langle \Sigma^{-1} \mu_0, \mu_0 \rangle - 2 \langle \Sigma^{-1} \mu, \mu_0 \rangle \right) \\ &= \frac{1}{2} \langle \Sigma^{-1}(\mu - \mu_0), \mu - \mu_0 \rangle. \end{aligned}$$

Der letzte Ausdruck ist gerade gleich  $\frac{1}{2} |\Sigma^{-1/2}(\mu - \mu_0)|^2$ . □

**Beispiel 3.5** Wir setzen Beispiel 3.1 unter der Annahme von i.i.d. Beobachtungsfehlern  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  fort. Der Vektor der Beobachtungen  $Y_i = f(x_i) + \varepsilon_i, i = 1, \dots, n$ , mit unbekannter Regressionsfunktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  ist dann gemäß  $\mathbb{P} = \mathcal{N}(\mu, \sigma^2 E_n)$  verteilt mit Mittelwertvektor  $\mu = (f(x_1), \dots, f(x_n))^T \in \mathbb{R}^n$ . In Modell  $p$  ist die Verteilungsfamilie gegeben durch  $(\mathbb{P}_{\beta^{(p)}})_{\beta^{(p)} \in \mathbb{R}^p}$  mit  $\mathbb{P}_{\beta^{(p)}} = \mathcal{N}(X^{(p)} \beta^{(p)}, \sigma^2 E_n)$  und der Designmatrix  $X^{(p)} = (e_j(x_i))_{i=1, \dots, n; j=1, \dots, p} \in \mathbb{R}^{n \times p}$ . Die Kullback-Leibler-Divergenz zwischen  $\mathbb{P}$  und  $\mathbb{P}_{\beta^{(p)}}$  ist damit gegeben durch

$$\begin{aligned} \text{KL}(\mathbb{P} | \mathbb{P}_{\beta^{(p)}}) &= \frac{1}{2\sigma^2} |\mu - X^{(p)} \beta^{(p)}|^2 \\ &= \frac{1}{2\sigma^2} \sum_{i=1}^n \left( f(x_i) - \sum_{j=1}^p \beta_j^{(p)} e_j(x_i) \right)^2 = \frac{n}{2\sigma^2} \left\| f - \sum_{j=1}^p \beta_j^{(p)} e_j \right\|_n^2. \end{aligned}$$

Der „Abstand“ zwischen  $\mathbb{P}$  und der bestmöglichen Wahl in  $(\mathbb{P}_{\beta^{(p)}})_{\beta^{(p)} \in \mathbb{R}^p}$  ist also klein, wenn  $f$  gut durch die ersten  $p$  Basisfunktionen approximiert werden kann. Als Kriterium zur Wahl von  $p$  können wir dies allerdings nicht direkt einsetzen, da  $f$  und damit auch  $\text{KL}(\mathbb{P} | \mathbb{P}_{\beta^{(p)}})$  unbekannt sind.

Für die Konvergenz in Lemma 3.3(iv) müssen wir nicht annehmen, dass das Wahrscheinlichkeitsmaß  $\mathbb{P}$ , unter dem die Beobachtungen erzeugt werden, in der Familie  $(\mathbb{P}_{\vartheta})_{\vartheta \in \Theta}$  enthalten ist. Wir benötigen lediglich  $\text{KL}(\mathbb{P} | \mathbb{P}_{\vartheta}) < \infty$  für  $\vartheta \in \Theta$ . Daher erhalten wir eine interessante Interpretation des *Maximum-Likelihood-Schätzers für misspezifizierte Modelle*: Während die linke Seite (die negative log-likelihood, skaliert mit  $1/n$ )

$$-\frac{1}{n} \sum_{i=1}^n l(\vartheta, X_i)$$

im Modell  $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$  durch den Maximum-Likelihood-Schätzer  $\widehat{\vartheta}$  minimiert wird, ist der Grenzwert

$$\text{KL}(\mathbb{P}|\mathbb{P}_\vartheta) - \text{KL}(\mathbb{P}|\mu)$$

unter demjenigen Parameter  $\vartheta$  kleinstmöglich, der die Kullback-Leibler-Divergenz  $\text{KL}(\mathbb{P}|\mathbb{P}_\vartheta)$  zum wahren  $\mathbb{P}$  minimiert. Dieser Wert ist damit ein natürlicher Kandidat für einen Grenzwert des Maximum-Likelihood-Schätzers. Zumindest für große Stichprobenumfänge können wir hoffen, dass  $\widehat{\vartheta}_p$  nahe am Minimierer von  $\text{KL}(\mathbb{P}|\mathbb{P}_\vartheta)$  liegt.

### 3.1.1 Akaike-Informationskriterium (AIC)

Es sei  $X \in \mathcal{X}$  eine Beobachtung aus dem Wahrscheinlichkeitsraum  $(\mathcal{X}, \mathcal{F}, \mathbb{P})$ . Wir betrachten dominierte statistische Modelle  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\vartheta^{(p)}})_{\vartheta^{(p)} \in \Theta_p})$  mit  $\Theta_p \subseteq \mathbb{R}^p$  mit nichtleerem Inneren (die Parametermenge  $\Theta_p$  ist  $p$ -dimensional).

Ohne Einschränkung können wir annehmen, dass alle Modelle durch dasselbe Maß  $\mu$  dominiert werden und auch  $\mathbb{P}$  absolutstetig bezüglich  $\mu$  ist. Im Modell  $(\mathbb{P}_{\vartheta^{(p)}})_{\vartheta^{(p)} \in \Theta_p}$  bezeichnen wir die Likelihood-Funktionen bezüglich  $\mu$  mit  $L_p(\vartheta^{(p)}, x)$ . Der Maximum-Likelihood-Schätzer  $\widehat{\vartheta}^{(p)}$  im  $p$ -ten Modell wird unter der Annahme  $X \sim \mathbb{P}_{\vartheta^{(p)}}$  für ein  $\vartheta^{(p)} \in \Theta_p$  bestimmt:

$$\widehat{\vartheta}^{(p)} \in \arg \max_{\vartheta^{(p)} \in \Theta_p} L_p(\vartheta^{(p)}, X).$$

Nachdem  $\mathbb{P}_{\widehat{\vartheta}^{(p)}}$  ein guter Kandidat innerhalb des  $p$ -ten Modells ist, wollen wir ein geeignetes Modell auswählen, indem wir die Kullback-Leibler-Divergenz  $\text{KL}(\mathbb{P}|\mathbb{P}_{\widehat{\vartheta}^{(p)}})$  über  $p$  minimieren. Da wir  $\mathbb{P}$  nicht kennen, ist jedoch  $\text{KL}(\mathbb{P}|\mathbb{P}_{\widehat{\vartheta}^{(p)}})$  nicht zugänglich und muss geschätzt werden. Als reelles Funktional wird  $\text{KL}(\mathbb{P}|\mathbb{P}_{\widehat{\vartheta}^{(p)}})$  dabei mit kleinerem Fehler schätzbar sein als  $\widehat{\vartheta}^{(p)}$  selbst. Hierfür verwenden wir im  $p$ -ten Modell für  $\vartheta^{(p)} \in \Theta_p$  die Darstellung (siehe Lemma 3.3)

$$\text{KL}(\mathbb{P}|\mathbb{P}_{\vartheta^{(p)}}) = \text{KL}(\mathbb{P}|\mu) - \mathbb{E}_{\mathbb{P}} \left[ \log L_p(\vartheta^{(p)}, X) \right],$$

sofern der letzte Erwartungswert existiert. In diesem Fall bezeichnen wir mit

$$d_p(\vartheta^{(p)}) := -2\mathbb{E}_{\mathbb{P}} [\log L_p(\vartheta^{(p)}, X)], \quad \vartheta^{(p)} \in \Theta_p,$$

die *Kullback-Leibler-Diskrepanz*.

Da  $\text{KL}(\mathbb{P}|\mu)$  unabhängig von  $\vartheta$  ist, ist die Minimierung von  $p \mapsto \text{KL}(\mathbb{P}|\mathbb{P}_{\widehat{\vartheta}^{(p)}})$  äquivalent zur Minimierung von  $p \mapsto d_p(\widehat{\vartheta}^{(p)})$ . Auch  $d_p(\vartheta)$  ist jedoch unbekannt.

Dessen empirische Version ist gegeben durch  $-2 \log L_p(\vartheta^{(p)}, X)$ , da  $X$  gemäß  $\mathbb{P}$  verteilt ist. Setzen wir in  $-\log L_p(\vartheta^{(p)}, X)$  den Maximum-Likelihood-Schätzer  $\widehat{\vartheta}^{(p)} = \widehat{\vartheta}^{(p)}(X)$  ein, erhalten wir aufgrund der Abhängigkeit zwischen  $\widehat{\vartheta}^{(p)}$  und  $X$  in der Regel keinen erwartungstreuen Schätzer von  $d_p(\vartheta^{(p)})$ . Stattdessen ist zu vermuten, dass  $d_p(\widehat{\vartheta}^{(p)})$  unterschätzt wird, denn  $\widehat{\vartheta}^{(p)}$  minimiert gerade die Funktion  $\vartheta^{(p)} \mapsto -\log L_p(\vartheta^{(p)}, X)$ .

*Beispiel 3.6 (Signalerkennung)* Für eine  $P$ -dimensionale mathematische Stichprobe  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} N(\vartheta, E_P)$  betrachten wir die Modelle

$$\Theta_p := \{\vartheta \in \mathbb{R}^P \mid \forall j > p : \vartheta_j = 0\}, \quad p = 1, \dots, P,$$

für den unbekannten Mittelwertvektor. Im  $p$ -ten Modell ist der Maximum-Likelihood-Schätzer durch das Stichprobenmittel in den ersten  $p$  Koordinaten gegeben, während die letzten  $P - p$  Koordinaten null sind:

$$\widehat{\vartheta}^{(p)} = (\widehat{\vartheta}_1, \dots, \widehat{\vartheta}_p, 0, \dots, 0) \quad \text{wobei} \quad \widehat{\vartheta}_j := (\overline{X}_n)_j = \frac{1}{n} \sum_{i=1}^n X_{i,j}$$

Schreiben wir  $|t|_{\leq p}^2 := \sum_{j=1}^p t_j^2$  und  $|t|_{> p}^2 := \sum_{j=p+1}^P t_j^2$  für  $t \in \mathbb{R}^P$ , dann gilt im  $p$ -ten Modell für  $\vartheta^{(p)} \in \Theta_p$

$$\begin{aligned} -2 \log L_p(\vartheta^{(p)}, X) &= -2 \sum_{i=1}^n \left( \log((2\pi)^{-\frac{P}{2}}) - \frac{|X_i - \vartheta^{(p)}|_{\leq p}^2 + |X_i|_{> p}^2}{2} \right) \\ &= \sum_{i=1}^n \left( |X_i - \vartheta^{(p)}|_{\leq p}^2 + |X_i|_{> p}^2 + P \log(2\pi) \right). \end{aligned}$$

Unter der Annahme  $X \sim \mathbb{P} = N(\vartheta_0, E_P)$  für ein  $\vartheta_0 \in \mathbb{R}^P$  erhalten wir mit der Bias-Varianz-Zerlegung

$$\begin{aligned} d_p(\vartheta^{(p)}) &= -2 \mathbb{E}_{\mathbb{P}}[\log L_p(\vartheta^{(p)}, X)] = \sum_{i=1}^n \left( \mathbb{E}_{\mathbb{P}}[|X_i - \vartheta^{(p)}|_{\leq p}^2 + |X_i|_{> p}^2] + P \log(2\pi) \right) \\ &= n \left( |\vartheta_0 - \vartheta^{(p)}|_{\leq p}^2 + |\vartheta_0|_{> p}^2 + P + P \log(2\pi) \right). \end{aligned}$$

Die Kullback-Leibler-Diskrepanz  $d_p$  spiegelt also wider, wie gut  $\vartheta_0$  durch einen Vektor aus  $\Theta_p$  approximiert werden kann: Auch wenn die Modellannahme  $\vartheta_0 \in \Theta_p$  nicht erfüllt ist, kann  $\Theta_p$  ein gutes Modell sein, wenn alle Einträge von  $\vartheta_0$  für große Indizes klein sind, das heißt wenn  $|\vartheta_0|_{> p}$  nahe 0 ist.

Setzen wir den Maximum-Likelihood-Schätzer  $\widehat{\vartheta}^{(p)}$  ein, erhalten wir

$$-2 \log L_p(\widehat{\vartheta}^{(p)}, X) = \sum_{i=1}^n \left( |X_i - \overline{X}_n|_{\leq p}^2 + |X_i|_{> p}^2 + P \log(2\pi) \right),$$

und aufgrund von  $\mathbb{E}[|X_i - \bar{X}|_{\leq p}^2] = p \frac{n-1}{n}$  gilt

$$\begin{aligned} -2\mathbb{E}_{\mathbb{P}}[\log L_p(\hat{\vartheta}^{(p)}, X)] &= n \left( p \frac{n-1}{n} + |\vartheta_0|_{>p}^2 + (P-p) + P \log(2\pi) \right) \\ &= n \left( -\frac{p}{n} + |\vartheta_0|_{>p}^2 + P + P \log(2\pi) \right). \end{aligned}$$

Andererseits gilt

$$\begin{aligned} \mathbb{E}_{\mathbb{P}}[d_p(\hat{\vartheta}^{(p)})] &= n \left( \mathbb{E}_{\mathbb{P}}[|\vartheta^0 - \bar{X}_n|_{\leq p}^2] + |\vartheta_0|_{>p}^2 + P + P \log(2\pi) \right) \\ &= n \left( \frac{p}{n} + |\vartheta_0|_{>p}^2 + P + P \log(2\pi) \right). \end{aligned}$$

Daraus folgt

$$\mathbb{E}_{\mathbb{P}}[-2 \log L_p(\hat{\vartheta}^{(p)}, X)] - \mathbb{E}_{\mathbb{P}}[d(\hat{\vartheta}^{(p)})] = -2p,$$

sodass  $-2 \log L_p(\hat{\vartheta}^{(p)}, X)$  die zu schätzende Größe  $d(\hat{\vartheta}^{(p)})$  systematisch um  $-2p$  unterschätzt.

Deutlich allgemeiner weisen wir den Zusammenhang  $\mathbb{E}_{\mathbb{P}}[-2 \log L(\hat{\vartheta}_p(X), X)] = \mathbb{E}_{\mathbb{P}}[d(\hat{\vartheta}_p)] - 2p$  in Satz 3.8 für das lineare Modell nach. Eine Biaskorrektur führt auf Akaikes Informationskriterium.

**Methode 3.7 (Modellwahl durch AIC)** Für die Modelle  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\vartheta})_{\vartheta \in \Theta_p})$  mit  $\Theta_p \subseteq \mathbb{R}^P$  und  $p = 1, \dots, P$  ist das **Akaike-Informationskriterium** (kurz: AIC) definiert als

$$\text{AIC}(p) := -2 \log L_p(\hat{\vartheta}^{(p)}, X) + 2p$$

für die Maximum-Likelihood-Schätzer  $\hat{\vartheta}^{(p)}$  im Modell  $(\mathbb{P}_{\vartheta})_{\vartheta \in \Theta_p}$ . Das Modell  $\hat{p}$  wird als Minimierer  $\hat{p} \in \arg \min_{p=1, \dots, P} \text{AIC}(p)$  gewählt und der resultierende, gemäß AIC ausgewählte Schätzer ist  $\hat{\vartheta}^{(\hat{p})}$ .

**Kurzbiografie (Hirotugu Akaike)** Hirotugu Akaike wurde 1927 in der Präfektur Shizuoka in Japan geboren. Er studierte an der Universität Tokio und promovierte dort 1961 am Institut für Statistische Mathematik. Akaike blieb auch in seiner weiteren wissenschaftlichen Laufbahn an diesem Institut, dessen Direktor er später wurde. Er schrieb wesentliche Arbeiten zur Theorie der Zeitreihen. Sein bekanntester Beitrag ist die Einführung des *Informationskriteriums* AIC im Jahr 1973, das später nach ihm benannt wurde. Akaike starb 2009.

Im Folgenden untersuchen wir, wie sich das Akaike-Informationskriterium auf das lineare Modell anwenden lässt.

**Satz 3.8 (AIC im Gaußschen linearen Modell)** Gegeben sei das wahre lineare Modell

$$Y = \mu + \varepsilon \quad \text{mit} \quad \mu \in \mathbb{R}^n, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 E_n)$$

unter  $\mathbb{P}$  mit bekannter Varianz  $\sigma^2 > 0$  sowie für  $p = 1, \dots, P$  mit  $P \leq n$  die Modelle

$$Y = X^{(p)}\beta^{(p)} + \varepsilon \quad \text{mit} \quad \beta^{(p)} \in \mathbb{R}^p, \quad \varepsilon \sim N(0, \sigma^2 E_n)$$

und die Designmatrix  $X^{(p)} \in \mathbb{R}^{n \times p}$  mit vollem Rang  $p$ . Dann gilt:

(i) Der Maximum-Likelihood-Schätzer in Modell  $p$  ist der Kleinste-Quadrate-Schätzer

$$\widehat{\beta}^{(p)} = (X^{(p)\top} X^{(p)})^{-1} X^{(p)\top} Y.$$

(ii) Das Akaike-Informationskriterium ist gegeben durch

$$\text{AIC}(p) = n \log(2\pi\sigma^2) + \frac{|Y - X^{(p)}\widehat{\beta}^{(p)}|^2}{\sigma^2} + 2p.$$

(iii)  $\text{AIC}(p)$  ist ein erwartungstreuer Schätzer der Kullback-Leibler-Diskrepanz:

$$\mathbb{E}_{\mathbb{P}}[\text{AIC}(p)] = \mathbb{E}_{\mathbb{P}}[d_p(\widehat{\beta}^{(p)})].$$

### Beweis

- (i) Da die Loglikelihood-Funktion von der Form  $\log L_p(\beta^{(p)}, Y) = \log((2\pi\sigma^2)^{-\frac{n}{2}}) - \frac{1}{2\sigma^2}|Y - X^{(p)}\beta^{(p)}|^2$  ist, entspricht der Maximum-Likelihood-Schätzer dem Kleinste-Quadrate-Schätzer, der in Lemma 2.14 explizit berechnet wurde.
- (ii) Einsetzen in die Definition von 3.7 liefert unmittelbar

$$\text{AIC}(p) = n \log(2\pi\sigma^2) + \sigma^{-2}|Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + 2p.$$

(iii) Wegen  $d_p(\beta^{(p)}) = \mathbb{E}[-2 \log L_p(\beta^{(p)}, Y)]$  folgt aus der Bias-Varianz-Zerlegung

$$\begin{aligned} d_p(\beta^{(p)}) &= n \log(2\pi\sigma^2) + \sigma^{-2} \mathbb{E}[|Y - X^{(p)}\beta^{(p)}|^2] \\ &= n \log(2\pi\sigma^2) + \sigma^{-2} (|\mu - X^{(p)}\beta^{(p)}|^2 + \underbrace{\mathbb{E}[|\varepsilon|^2]}_{=\sigma^2 n}) \\ &= n(\log(2\pi\sigma^2) + 1) + \sigma^{-2}|\mu - X^{(p)}\beta^{(p)}|^2. \end{aligned}$$

Einsetzen von  $X^{(p)}\widehat{\beta}^{(p)} = \Pi^{(p)}Y = \Pi^{(p)}\mu + \Pi^{(p)}\varepsilon$  mit der Orthogonalprojektion  $\Pi^{(p)} := \Pi_{X^{(p)}}$  liefert

$$\begin{aligned} \mathbb{E}[d_p(\widehat{\beta}^{(p)})] &= n(\log(2\pi\sigma^2) + 1) + \sigma^{-2} \mathbb{E}[|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2] \\ &= n(\log(2\pi\sigma^2) + 1) + \sigma^{-2} (|\mu - \Pi^{(p)}\mu|^2 + \mathbb{E}[|\Pi^{(p)}\varepsilon|^2]) \\ &= n(\log(2\pi\sigma^2) + 1) + \sigma^{-2} |(E_n - \Pi^{(p)})\mu|^2 + p, \end{aligned}$$

da  $\Pi^{(p)}$  Rang  $p$  hat und damit  $\mathbb{E}[|\Pi^{(p)}\varepsilon|^2] = \sigma^2 p$  gilt. Analog ist  $E_n - \Pi^{(p)}$  die orthogonale Projektion auf einen  $(n-p)$ -dimensionalen Unterraum und  $\mathbb{E}[|(E_n - \Pi^{(p)})\varepsilon|^2] = (n-p)\sigma^2$ . Wir erhalten

$$\begin{aligned}
\mathbb{E}[\text{AIC}(p)] &= n \log(2\pi\sigma^2) + \sigma^{-2} \mathbb{E}[|Y - X^{(p)}\widehat{\beta}^{(p)}|^2] + 2p \\
&= n \log(2\pi\sigma^2) + \sigma^{-2} (|\mu - \Pi^{(p)}\mu|^2 + \mathbb{E}[|\varepsilon - \Pi^{(p)}\varepsilon|^2]) + 2p \\
&= \mathbb{E}[d_p(\widehat{\beta}^{(p)})].
\end{aligned}$$

Damit ist die Erwartungstreue von  $\text{AIC}(p)$  gezeigt.  $\square$

Die Darstellung von  $\text{AIC}(p)$  erinnert an die Bias-Varianz-Zerlegung. Durch Multiplikation mit  $\sigma^2$  (was für die Minimierung von  $\text{AIC}(p)$  nichts ändert) erhalten wir die quadrierten Residuen  $|Y - X^{(p)}\widehat{\beta}^{(p)}|^2$  als Maß für die Modellgüte sowie den doppelten Varianzterm  $2p\sigma^2$ . Die Modellwahl ist gegeben durch

$$\begin{aligned}
\widehat{p} &\in \arg \min_p \left( n \log(2\pi\sigma^2) + \frac{|Y - X^{(p)}\widehat{\beta}^{(p)}|^2}{\sigma^2} + 2p \right) \\
&= \arg \min_p \underbrace{(|Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + 2p\sigma^2)}_{\text{AIC}_{\text{LM}}(p)}.
\end{aligned}$$

Hierbei fällt der empirische Verlust  $|Y - X^{(p)}\widehat{\beta}^{(p)}|^2$  in  $\text{AIC}_{\text{LM}}(p)$  mit wachsender Dimension  $p$ , während der Strafterm  $2p\sigma^2$  in  $p$  wächst.  $\widehat{p}$  balanciert also die Güte der Datenanpassung mit der Komplexität des Modells aus. Man beachte, dass die Erwartungstreue von  $\text{AIC}(p)$  allein nicht sicherstellt, dass  $\widehat{p}$  tatsächlich ein gutes Modell bzw.  $X^{(\widehat{p})}\widehat{\beta}^{(\widehat{p})}$  eine gute Vorhersage für  $\mu$  ist. Die Qualität von  $X^{(\widehat{p})}\widehat{\beta}^{(\widehat{p})}$  beurteilen wir in Kapitel 3.2 anhand einer Orakelungleichung.

**Korollar 3.9 (Unverzerrte Risikoschätzung)** *In der Situation von Satz 3.8 ist  $|Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + 2p\sigma^2 - n\sigma^2$  eine erwartungstreue Schätzung des quadratischen Vorhersagefehlers  $|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2$ :*

$$\mathbb{E}_{\mathbb{P}}[|Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + 2p\sigma^2 - n\sigma^2] = \mathbb{E}_{\mathbb{P}}[|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2]$$

*Insbesondere kann das Akaike-Informationskriterium als „unbiased risk estimation“-Kriterium interpretiert werden.*

**Beweis** Es gilt für die quadrierten Residuen wie oben gesehen, dass

$$\mathbb{E}_{\mathbb{P}}[|Y - X^{(p)}\widehat{\beta}^{(p)}|^2] = |\mu - \Pi^{(p)}\mu|^2 + (n - p)\sigma^2.$$

Mit der Bias-Varianz-Zerlegung folgt

$$\mathbb{E}_{\mathbb{P}}[|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2] = |\mu - \Pi^{(p)}\mu|^2 + p\sigma^2.$$

Ein Vergleich der beiden Erwartungswerte liefert die Behauptung.  $\square$

**Bemerkung 3.10 (Mallows'  $C_p$ -Kriterium)** Der Beweis von Korollar 3.9 zeigt, dass für die erwartungstreue Fehlerschätzung die Normalverteilungsannahme unerheblich ist. Es genügt,  $\mathbb{E}[\varepsilon] = 0$  und  $\text{Cov}(\varepsilon) = \sigma^2 E_n$  vorauszusetzen. Das Modell allgemein

aufgrund einer erwartungstreuen Risikoschätzung auszuwählen, ist der Ansatz von Mallows'  $C_p$ -Kriterium, wobei

$$C_p := |Y - X^{(p)} \widehat{\beta}^{(p)}|^2 + 2p\sigma^2 - n\sigma^2$$

gilt. Bis auf einen von  $p$  unabhängigen Summanden stimmt Mallows'  $C_p$ -Kriterium im linearen Modell mit AIC überein.

Wir wollen nun Satz 3.8 auf den Fall einer unbekannten Varianz erweitern:

**Satz 3.11 (AIC im linearen Modell mit unbekannter Varianz)** *Gegeben sei das wahre lineare Modell*

$$Y = \mu + \varepsilon \quad \text{mit} \quad \mu \in \mathbb{R}^n, \quad \varepsilon \sim N(0, \sigma_0^2 E_n)$$

unter  $\mathbb{P}$  sowie für  $p = 1, \dots, P$  mit  $P \leq n - 3$  die Modelle

$$Y = X^{(p)} \beta^{(p)} + \varepsilon \quad \text{mit} \quad \beta^{(p)} \in \mathbb{R}^p, \quad \varepsilon \sim N(0, \sigma^2 E_n)$$

und der Designmatrix  $X^{(p)} \in \mathbb{R}^{n \times p}$  mit vollem Rang  $p$ . Für  $k = p + 1$  ist  $\vartheta^{(k)} = (\beta^{(p)}, \sigma^2)$  der unbekannte Parameter in  $\mathbb{R}^k$ . Dann gilt:

(i) Der Maximum-Likelihood-Schätzer ist  $\widehat{\vartheta}^{(k)} = (\widehat{\beta}^{(p)}, \widehat{\sigma}_p^2)$  mit dem Kleinsten-Quadrate-Schätzer  $\widehat{\beta}^{(p)} \in \mathbb{R}^p$  und  $\widehat{\sigma}_k^2 := \frac{1}{n} |Y - X^{(p)} \widehat{\beta}^{(p)}|^2$ . Es gilt

$$\text{AIC}(k) = n \log(2\pi \widehat{\sigma}_k^2) + n + 2k.$$

(ii) Im Fall, dass  $\mu = X^{(p)} \beta_0^{(p)}$  für ein  $p$  und ein  $\beta_0^{(p)} \in \mathbb{R}^p$  erfüllt ist, gilt mit  $k = p + 1$

$$\mathbb{E}[\text{AIC}(k)] = \mathbb{E}[d_k(\widehat{\vartheta}_k)] - 2 \frac{k(k+1)}{n-k-1}.$$

Für eine unbekannte Varianz liefert  $\text{AIC}(k)$  also keine unverzerrte Schätzung der Kullback-Leibler-Diskrepanz. Für allgemeines  $\mu \in \mathbb{R}^n$  in (ii) kann man  $\mathbb{E}[d_k(\widehat{\vartheta}_k)] - \mathbb{E}[\text{AIC}(k)] \in [0, 2k(k+1)/(n-k-1)]$  zeigen. Die Abweichung ist also von der Größenordnung  $O(k^2/n)$  in der Parameterdimension  $k$  und der Stichprobengröße  $n$ , was für große  $n$  bei beschränkter maximaler Modellgröße  $P$  vernachlässigbar ist. Ähnliches gilt für AIC in allgemeinen Likelihood-Modellen.

### **Beweis**

(i) Die Loglikelihood-Funktion erfüllt

$$-2 \log L_k(\vartheta^{(k)}) = n \log((2\pi\sigma^2)) + \frac{1}{\sigma^2} |Y - X^{(p)} \beta^{(p)}|^2,$$

sodass die Form des Maximum-Likelihood-Schätzers klar ist, siehe auch Aufgabe 2.6. Einsetzen liefert

$$\begin{aligned} \text{AIC}(k) &= -2 \log L_k(\widehat{\vartheta}^{(k)}) + 2k = n \log(2\pi\widehat{\sigma}_k^2) + \frac{1}{\widehat{\sigma}_k^2} |Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + 2k \\ &= n \log(2\pi\widehat{\sigma}_k^2) + n + 2k. \end{aligned}$$

(ii) Für  $\vartheta^{(k)} = (\beta^{(p)}, \sigma^2)$  gilt

$$d_k(\vartheta^{(k)}) = \mathbb{E}_{\mathbb{P}}[-2 \log L_k(\beta^{(p)}, \sigma^2)] = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} (|\mu - X^{(p)}\beta^{(p)}|^2 + n\sigma_0^2),$$

sodass

$$\mathbb{E}_{\mathbb{P}}[d_k(\widehat{\vartheta}^{(k)})] = n\mathbb{E}[\log(2\pi\widehat{\sigma}_k^2)] + \mathbb{E}\left[\frac{1}{\widehat{\sigma}_k^2} (|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2 + n\sigma_0^2)\right].$$

Wegen  $\mathbb{E}_{\mathbb{P}}[\text{AIC}(k)] = n\mathbb{E}[\log(2\pi\widehat{\sigma}_k^2)] + n + 2k$  folgt

$$\begin{aligned} \mathbb{E}_{\mathbb{P}}[\text{AIC}(k)] - \mathbb{E}_{\mathbb{P}}[d_k(\widehat{\vartheta}^{(k)})] &= n + 2k - \mathbb{E}\left[\frac{1}{\widehat{\sigma}_k^2} (|\mu - X^{(p)}\widehat{\beta}^{(p)}|^2 + n\sigma_0^2)\right] \\ &= n + 2k - \mathbb{E}\left[\frac{|\mu - \Pi^{(p)}Y|^2 + n\sigma_0^2}{\frac{1}{n}|(E_n - \Pi^{(p)})Y|^2}\right]. \end{aligned}$$

Da  $\Pi^{(p)}$  und  $E_n - \Pi^{(p)}$  Orthogonalprojektion auf  $\text{Im } X^{(p)}$  bzw.  $\text{Im}(X^{(p)})^\perp$  sind, sind  $\Pi^{(p)}Y$  und  $(E_n - \Pi^{(p)})Y$  unkorreliert und damit unabhängig. Es folgt

$$\begin{aligned} \mathbb{E}_{\mathbb{P}}[\text{AIC}(k)] - \mathbb{E}_{\mathbb{P}}[d_k(\widehat{\vartheta}^{(k)})] &= n + 2k - \mathbb{E}[|\mu - \Pi^{(p)}Y|^2 + n\sigma_0^2] \mathbb{E}\left[\frac{n}{|(E_n - \Pi^{(p)})Y|^2}\right] \\ &= n + 2k - (|\mu - \Pi^{(p)}\mu|^2 + p\sigma_0^2 + n\sigma_0^2) \frac{n}{\sigma_0^2} \mathbb{E}[Z^{-1}], \end{aligned}$$

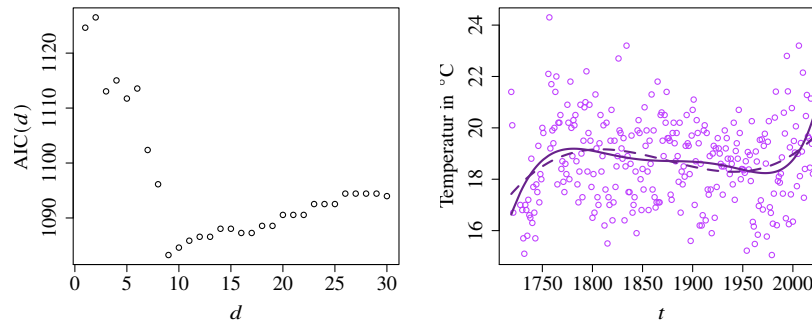
wobei  $Z := \frac{1}{\sigma_0^2} |(E_n - \Pi^{(p)})Y|^2 = \frac{1}{\sigma_0^2} |(E_n - \Pi^{(p)})\varepsilon|^2 \sim \chi^2(n-p)$  und  $|\mu - \Pi^{(p)}\mu|^2 = 0$  im korrekt spezifizierten Modell gilt. Wir berechnen

$$\mathbb{E}[Z^{-1}] = \int_0^\infty \frac{1}{z} \frac{2^{-\frac{n-p}{2}}}{\Gamma(\frac{n-p}{2})} z^{\frac{n-p}{2}-1} e^{-\frac{z}{2}} dz = \frac{1}{n-p-2} = \frac{1}{n-k-1}$$

für  $n \geq p+3$  und erhalten

$$\mathbb{E}_{\mathbb{P}}[\text{AIC}(k)] - \mathbb{E}_{\mathbb{P}}[d_k(\widehat{\vartheta}^{(k)})] = n + 2k - \frac{(k-1+n)n}{n-k-1} = -2 \frac{k(k+1)}{n-k-1}.$$

*Beispiel 3.12 (AIC für Klimadaten)* Wir erinnern uns an die Klimadatenzeitreihe aus Beispiel 2.52, in der wir die mittleren Julitemperaturen in Berlin-Dahlem zwischen 1719 und 2020 betrachtet haben. Wir hatten hier Regressionspolynome mit den Graden  $d \in \{1, \dots, 4\}$  verwendet.  $\text{AIC}(d)$  für  $d \in \{1, \dots, 30\}$  ist in Abbildung 3.1



**Abb. 3.1** Links: AIC für polynomielle Regression in Abhängigkeit vom Polynomgrad für die Temperaturdaten aus Beispiel 2.52. Rechts: Das Regressionpolynom des Grades 9 (durchgezogene Linie) sowie des Grades 3 (gestrichelte Linie). Datenquelle: Deutscher Wetterdienst

dargestellt. Bis  $d = 9$  fällt das Informationskriterium im Wesentlichen, hier profitieren wir also durch die Erhöhung der Parameterdimension noch deutlich in den Residuen. Ab  $d = 10$  wächst  $AIC(d)$ , sodass ab hier der Einfluss des Strafterms dominiert. Ebenfalls ist in Abbildung 3.1 das durch AIC gewählte Regressionspolynom zusammen mit dem Polynom vom Grad 3, für das wir uns in Beispiel 2.52 entschieden hatten, zu sehen.

### 3.1.2 Das Bayes-Informationskriterium (BIC)

Anstelle des Ansatzes, die Kullback-Leibler-Diskrepanz oder das Risiko erwartungstreu zu schätzen, ist der Bayes-Ansatz (der auf Gideon Schwarz zurückgeht), eine gleichmäßige a-priori-Verteilung für die Modelle  $p \in \{1, \dots, P\}$  anzunehmen. Die a-posteriori-Verteilung liefert uns dann ein Kriterium zur Auswahl des Modells  $p$ .

**Kurzbiografie (Gideon Schwarz)** Gideon E. Schwarz wurde 1933 in Salzburg geboren. Er flüchtete nach dem Anschluss Österreichs 1938 in das damaligen Palästina, heute Israel. Gideon E. Schwarz studierte an der Hebrew University Mathematik. Seinen Ph.D. erwarb er an der Columbia University 1961 im Bereich der mathematischen Statistik. Später wurde er Professor an der Hebrew University mit zahlreichen Gastaufenthalten in Stanford, Tel Aviv und Berkeley. Er verfasste vor allem wichtige Beiträge zur Bayes-Statistik. Als Reaktion auf das *Informationskriterium* AIC entwickelte er 1976 das *Bayes-Informationskriterium* BIC, manchmal auch Schwarz Information Criterion (SIC) genannt. Schon damals entwickelte sich eine große Diskussion über die verschiedenen Sichtweisen, die teilweise bis heute anhält. Schwarz starb im Jahr 2007.

Um das Bayes-Informationskriterium (englisch: *Bayesian information criterion*, kurz: BIC) herzuleiten, betrachten wir wieder das wahre lineare Modell

$$Y = \mu + \varepsilon \quad \text{mit} \quad \varepsilon \sim N(0, \sigma^2 E_n)$$

und dem unbekannten Mittelwertvektor  $\mu \in \mathbb{R}^n$ , den es zu schätzen gilt. Das Rauschniveau  $\sigma > 0$  nehmen wir als bekannt an. Unser Ziel ist, aus den Modellen

$$Y = X^{(p)} \beta^{(p)} + \varepsilon$$

mit  $\beta^{(p)} \in \mathbb{R}^p$ ,  $X^{(p)} \in \mathbb{R}^{n \times p}$ ,  $\text{rank } X = p$  und  $\varepsilon \sim N(0, \sigma^2 E_n)$  für  $p \in \{1, \dots, P\}$  ein geeignetes auszuwählen.

Für jedes  $p$  sei die a-priori-Verteilung des Parameters  $\beta^{(p)}$  durch eine Lebesgue-Dichte  $\pi_p$  auf  $\mathbb{R}^p$  gegeben. Für das Modell  $p$  wählen wir die uninformativ a-priori-Verteilung  $U := U(\{1, \dots, P\})$ , sodass jedes Modell  $p$  mit der gleichen Wahrscheinlichkeit  $\frac{1}{P}$  ausgewählt wird. Das gemäß  $U$  ausgewählte Modell bezeichnen wir mit  $\kappa$ . Die a-priori-Verteilungen von  $p$  und  $\beta^{(p)}$  seien unabhängig.

Die gemeinsame Verteilung von Modell  $\kappa$ , Parameter  $\beta$  und Beobachtung  $Y$  ist dann gegeben durch

$$\mathbb{P}^{(\kappa, \beta, Y)}(A) = \int \mathbb{1}_A(p, \beta^{(p)}, y) \frac{e^{-|y - X^{(p)} \beta^{(p)}|^2 / (2\sigma^2)}}{(2\pi\sigma^2)^{n/2}} dy \pi_p(\beta^{(p)}) d\beta^{(p)} U(dp).$$

Da die Dimension von  $\beta^{(p)}$  von  $p$  abhängt, ist in dieser Formel Vorsicht geboten. Um einen gemeinsamen Grundraum zu haben, müssen wir  $A \subseteq \{1, \dots, P\} \times \mathbb{R}^P \times \mathbb{R}^n$  wählen und in der Indikatorfunktion  $\beta^{(p)}$  als Vektor in  $\mathbb{R}^P$  interpretieren, dessen letzte  $P - p$  Koordinaten null sind. Da wir nur an der gemeinsamen Verteilung von  $\kappa$  und  $Y$  interessiert sind, entfällt dieses technische Problem ohnehin: Die Verteilung von  $(\kappa, Y)$  besitzt die Dichte

$$f^{(\kappa, Y)}(p, y) = \int_{\mathbb{R}^p} \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{|y - X^{(p)} \beta^{(p)}|^2}{2\sigma^2}\right) \pi_p(\beta^{(p)}) d\beta^{(p)}$$

bezüglich des Produktmaßes  $U \otimes \lambda^{\otimes n}$ . Der BIC-Ansatz besteht nun darin, das Modell mit der größten a-posteriori-Wahrscheinlichkeit auszuwählen, siehe MAP-Schätzer aus Korollar 1.33:

$$\hat{p} \in \arg \min_{p=1, \dots, P} \mathbb{P}(\kappa = p | Y)$$

Für dieses Kriterium leiten wir nun eine deutlich zugänglichere (approximative) Darstellung her. Für die a-posteriori-Verteilung von  $\kappa$  gilt

$$\mathbb{P}(\kappa = p | Y) \propto \int_{\mathbb{R}^k} \exp\left(-\frac{|Y - X^{(p)} \beta^{(p)}|^2}{2\sigma^2}\right) \pi_p(\beta^{(p)}) d\beta^{(p)}.$$

Es bezeichne wieder  $\hat{\beta}^{(p)}$  den Kleinste-Quadrate-Schätzer in Modell  $p$ . Aufgrund der Orthogonalität der Residuen  $Y - X^{(p)} \hat{\beta}^{(p)} = (E_n - \Pi_{X^{(p)}})Y$  zum Bild von  $X^{(p)}$  können wir den Exponenten in der a-posteriori-Verteilung gemäß

$$|Y - X^{(p)} \beta^{(p)}|^2 = |Y - X^{(p)} \hat{\beta}^{(p)}|^2 + |X^{(p)} \hat{\beta}^{(p)} - X^{(p)} \beta^{(p)}|^2$$

$$\begin{aligned}
& -2\langle Y - X^{(p)}\beta^{(p)}, X^{(p)}(\widehat{\beta}^{(p)} - \beta^{(p)}) \rangle \\
& = |Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + |X^{(p)}\widehat{\beta}^{(p)} - X^{(p)}\beta^{(p)}|^2
\end{aligned}$$

zerlegen. Da  $\widehat{\beta}^{(p)}$  nur von  $Y$  (und  $p$ ), aber nicht von  $\beta$  abhängt, erhalten wir

$$\begin{aligned}
\mathbb{P}(\kappa = p|Y) & \propto \exp\left(-\frac{1}{2\sigma^2}|Y - X^{(p)}\widehat{\beta}^{(p)}|^2\right) \\
& \times \int_{\mathbb{R}^p} \exp\left(-\frac{1}{2\sigma^2}|X^{(p)}(\widehat{\beta}^{(p)} - \beta^{(p)})|^2\right) \pi_p(\beta^{(p)}) d\beta^{(p)}.
\end{aligned}$$

Durch eine Substitution  $h = \sqrt{n}(\widehat{\beta}^{(p)} - \beta^{(p)})$  erhalten wir für  $\Sigma_n^{(p)} := \frac{1}{n}(X^{(p)})^\top X^{(p)} \in \mathbb{R}^{p \times p}$

$$\begin{aligned}
& \int_{\mathbb{R}^p} \exp\left(-\frac{1}{2\sigma^2}|X^{(p)}(\widehat{\beta}^{(p)} - \beta^{(p)})|^2\right) \pi_p(\beta^{(p)}) d\beta^{(p)} \\
& = n^{-p/2} \int_{\mathbb{R}^p} \exp\left(-\frac{1}{2\sigma^2}\langle \Sigma_n^{(p)} h, h \rangle\right) \pi_p\left(\widehat{\beta}^{(p)} - \frac{h}{\sqrt{n}}\right) dh =: n^{-p/2} I_{n,p}.
\end{aligned}$$

Bezeichnen wir die Normierungskonstante der Dichte von  $\mathbb{P}(\kappa = p|Y)$  mit  $C > 0$ , so ergibt sich

$$\log \mathbb{P}(\kappa = p|Y) = -\log C - \frac{p}{2} \log n - \frac{1}{2\sigma^2} |Y - X^{(p)}\widehat{\beta}^{(p)}|^2 + \log I_{n,p} \quad (3.2)$$

Um den letzten Term  $\log I_{n,p}$  abzuschätzen, [RD]betrachten wir die Asymptotik  $n \rightarrow \infty$  und nehmen an, dass

$$\Sigma_n^{(p)} \rightarrow \Sigma^{(p)} \quad \text{für} \quad n \rightarrow \infty$$

konvergiert.

*Beispiel 3.13 (Asymptotik der Designmatrix)* Es seien  $\varphi_1, \dots, \varphi_n: [0, 1] \rightarrow \mathbb{R}$  stetige Funktionen, beispielsweise die Monome  $\varphi_k(x) = x^{k-1}$ . Unter äquidistantem Design ist die Designmatrix durch

$$X^{(p)} = \begin{pmatrix} \varphi_1(1/n) & \cdots & \varphi_p(1/n) \\ \varphi_1(2/n) & \cdots & \varphi_p(2/n) \\ \vdots & & \vdots \\ \varphi_1(n/n) & \cdots & \varphi_p(n/n) \end{pmatrix}$$

gegeben, und wir erhalten für alle  $i, j \in \{1, \dots, p\}$

$$(\Sigma_n^{(p)})_{i,j} = \frac{1}{n} \sum_{k=1}^n \varphi_i\left(\frac{k}{n}\right) \varphi_j\left(\frac{k}{n}\right) \rightarrow \int_0^1 \varphi_i(x) \varphi_j(x) dx =: (\Sigma^{(p)})_{i,j}$$

für  $n \rightarrow \infty$ .

Sind nun noch die Dichten  $\pi_p$  gleichmäßig beschränkt, finden wir mit dominierter Konvergenz eine obere Schranke für  $I_{n,p}$ , die gleichmäßig in  $n$  und  $p$  gilt. Aus (3.2) folgt dann für eine Nullfolge  $o(1) \rightarrow 0$  für  $n \rightarrow \infty$

$$-2 \log \mathbb{P}(\kappa = p | Y) = p \log n(1 + o(1)) + \frac{1}{\sigma^2} |Y - X^{(p)} \widehat{\beta}^{(p)}|^2.$$

Vernachlässigt man den  $o(1)$ -Term, entspricht das Bayes-Informationskriterium gerade

$$\widehat{p} \in \arg \min_{p=1, \dots, P} \left\{ \frac{1}{\sigma^2} |Y - X^{(p)} \widehat{\beta}^{(p)}|^2 + p \log n \right\}.$$

Die quadrierten Residuen entsprechen wieder der negativen Loglikelihood-Funktion im linearen Modell (mit Faktor  $2\sigma^2$ ). Hinzu kommt eine Penalisierung mit  $p \log n$ . Analog zu Akaikes Informationskriterium definieren wir daher ganz allgemein:

**Methode 3.14 (Modellwahl durch BIC)** Für  $p = 1, \dots, P$  und dominierte statistische Modelle  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\theta^{(p)}})_{\theta^{(p)} \in \Theta_p})$  mit  $\mathcal{X} = \mathbb{R}^n$ ,  $\mathcal{F} = \mathcal{B}(\mathbb{R}^n)$  und  $\Theta_p \subseteq \mathbb{R}^p$  ist das **Bayes-Informationskriterium** (kurz: BIC) definiert als

$$\text{BIC}(p) := -2 \log L_p(\widehat{\theta}^{(p)}, X) + p \log n$$

für die Likelihood-Funktion  $L_p$  und den Maximum-Likelihood-Schätzer  $\widehat{\theta}_p$  im Modell  $p$ . Das Modell  $\widehat{p}$  wird als Minimierer  $\widehat{p} \in \arg \min_{p=1, \dots, P} \text{BIC}(p)$  gewählt.

Entsprechend der vorangegangenen Herleitung erhalten wir:

**Lemma 3.15** Für  $p = 1, \dots, P$  mit  $P \leq n$  und die linearen Modelle

$$Y = X^{(p)} \beta^{(p)} + \varepsilon \quad \text{mit} \quad \beta^{(p)} \in \mathbb{R}^p, \quad \varepsilon \sim N(0, \sigma^2 E_n),$$

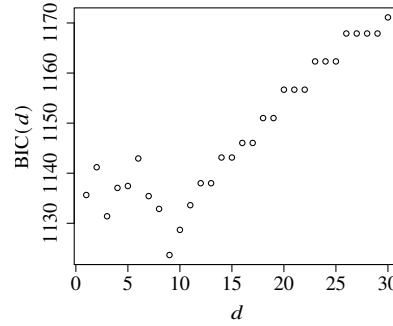
Designmatrix  $X^{(p)} \in \mathbb{R}^{n \times p}$  von vollem Rang  $p$  und mit bekannter Varianz  $\sigma > 0$  gilt

$$\text{BIC}(p) = n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} |Y - X^{(p)} \widehat{\beta}^{(p)}|^2 + p \log n.$$

Im Vergleich zum AIC wird im BIC die Modelldimension mit  $p \log n$  statt mit  $2p$  bestraft. BIC wählt also für  $n \geq 8$  tendenziell ein kleineres Modell aus. Andererseits führt die AIC-Wahl typischerweise auf einen kleineren Vorhersagefehler  $|Y - X^{(p)} \widehat{\beta}^{(p)}|^2$ .

Man beachte, dass  $|Y - X^{(p)} \widehat{\beta}^{(p)}|^2$  in  $n$  wächst und ohne weitere Annahmen von der Größenordnung  $n$  ist. Reskalieren wir BIC durch  $\frac{1}{n} |Y - X^{(p)} \widehat{\beta}^{(p)}|^2 + p \frac{\log n}{n}$ , wird deutlich, dass die Erhöhung der Modelldimension mit steigendem  $n$  geringer bestraft wird.

**Beispiel 3.16 (BIC für Klimadaten)** Wir wenden das Bayes-Informationskriterium auf die Daten aus den Beispielen 2.52 bzw. 3.12 über die langfristige Entwicklung



**Abb. 3.2** BIC für polynomielle Regression in Abhängigkeit vom Polynomgrad für die Temperaturdaten aus Beispiel 2.52

der mittleren Julitemperaturen an. Die resultierenden Werte  $BIC(d)$  für die Polynomgrade  $d \in \{1, \dots, 30\}$  sind in Abbildung 3.2 dargestellt. Im Vergleich zu den AIC-Werten steigt BIC wesentlich schneller in  $d$  an. Zwar wählt auch das Bayes-Informationskriterium den Grad 9, aber der Abstand zum Grad 3 ist hier viel kleiner als bei Akaiikes Informationskriterium.

### 3.2 Orakelungleichung für die penalisierte Modellwahl

Im vorangegangenen Abschnitt haben wir die Informationskriterien AIC und BIC hergeleitet und Beispiele gesehen, in denen diese gute Resultate erzielen. Die offene Frage ist, ob wir im Allgemeinen nachweisen können, dass die Modellwahlkriterien den Trade-off zwischen Datenanpassung und zu hoher Modellkomplexität lösen.

Wie zuvor seien die Beobachtungen durch

$$Y = \mu + \varepsilon, \quad \mu \in \mathbb{R}^n, \varepsilon \sim N(0, \sigma^2 E_n),$$

gegeben. Für die Modelle  $Y = X^{(p)}\beta^{(p)} + \varepsilon$  mit den Designmatrizen  $X^{(p)} \in \mathbb{R}^{n \times p}$  von vollem Rang bezeichne wieder  $\hat{\beta}^{(p)}$  den jeweiligen Kleinste-Quadrate-Schätzer, wobei  $p \in \{1, \dots, P\}$ . Die Modellwahlkriterien AIC und BIC zur Wahl des finalen Schätzers  $\hat{\beta}^{(\hat{p})}$  waren beide von der Form

$$\hat{p} \in \arg \min_{p=1, \dots, P} \left( \|Y - X^{(p)}\hat{\beta}^{(p)}\|^2 + \text{Pen}(p) \right)$$

mit dem Penalisierungsterm  $\text{Pen}(p)$ . Für AIC ist Letzterer durch  $\text{Pen}(p) = 2\sigma^2 p$  gegeben, während BIC dem Strafterm  $\text{Pen}(p) = \log(n)\sigma^2 p$  entspricht. In beiden Fällen hängt der Strafterm also von der Dimension des Parameterraums ab. Betrachtet man

$$S_p := \mathfrak{I}X^{(p)} \quad \text{und} \quad \hat{\mu}^{(p)} := X^{(p)}\hat{\beta}^{(p)} = \Pi_{S_p} Y$$

mit der Orthogonalprojektion  $\Pi_{S_p}$  auf das Bild von  $X^{(p)}$ , betten sich beide Verfahren in die folgende allgemeine Modellwahl ein.

**Definition 3.17** Für eine Beobachtung  $Y \in \mathbb{R}^n$  seien lineare Unterräume  $S_m \subseteq \mathbb{R}^n$  der Dimension  $p_m$  gegeben, wobei  $m = 1, \dots, M$ . Für eine monoton wachsende Funktion  $\text{Pen}: \mathbb{N} \rightarrow [0, \infty)$  betrachten wir die **penalisierte Modellwahl**

$$\hat{m} \in \arg \min_{m=1, \dots, M} \left( |Y - \hat{\mu}^{(m)}|^2 + \text{Pen}(p_m) \right)$$

mit den Kleinst-Quadrate-Schätzern  $\hat{\mu}^{(m)} := \Pi_{S_m} Y \in S_m$ . Wir setzen  $\mu^{(m)} = \Pi_{S_m} \mu$ .

Wir werden nun eine Orakelungleichung für die penalisierte Modellwahl beweisen. Eine erste Orakelungleichung ist uns bereits in Kapitel ?? begegnet. Mit ihrer Hilfe konnten wir das Verhalten des Kleinst-Quadrate-Schätzers im misspezifizierten Modell analysieren, siehe Korollar ?. Dieses Resultat müssen wir um den Penalisierungsterm erweitern. Die größte Schwierigkeit besteht darin, die Abhängigkeit zwischen  $\hat{m}$  und  $\hat{\mu}^{(m)}$  im gewählten Parameter  $\hat{\mu}^{(\hat{m})}$  zu kontrollieren.

Ausgangspunkt für den Beweis der Orakelungleichung war die Fundamentalungleichung aus Lemma (??). Im penalisierten Fall (und für quadratischen Verlust) erhalten wir folgendes Analogon:

**Lemma 3.18 (Fundamentalungleichung)** *Im Modell  $Y = \mu + \varepsilon$  gilt für den penalisierten empirischen Risikominimierer*

$$\hat{m} := \arg \min_{m=1, \dots, M} \left\{ |Y - \hat{\mu}^{(m)}|^2 + \text{Pen}(\bar{p}_m) \right\}$$

die Ungleichung

$$|\mu - \hat{\mu}^{(\hat{m})}|^2 + \text{Pen}(p_{\hat{m}}) \leq |\mu - \mu^{(m)}|^2 + \text{Pen}(p_m) + 2\langle \varepsilon, \hat{\mu}^{(\hat{m})} - \mu^{(m)} \rangle$$

für jedes  $m = 1, \dots, M$ .

**Beweis** Nach Definition von  $\hat{m}$  gilt für jedes  $m$

$$|Y - \hat{\mu}^{(\hat{m})}|^2 + \text{Pen}(p_{\hat{m}}) \leq |Y - \hat{\mu}^{(m)}|^2 + \text{Pen}(p_m) \leq |Y - \mu^{(m)}|^2 + \text{Pen}(p_m),$$

wobei wir  $|Y - \hat{\mu}^{(m)}|^2 = |Y - \Pi_{S_m} Y|^2 \leq |Y - \mu^{(m)}|^2$  aus der Eigenschaft der Orthogonalprojektion  $\Pi_{S_m}$  und  $\mu^{(m)} \in S_m$  folgern. Aus  $|Y - v|^2 = |v - \mu|^2 + |\varepsilon|^2 - 2\langle \varepsilon, v - \mu \rangle$  für jedes  $v \in \mathbb{R}^n$  folgt durch Einsetzen

$$|\mu - \hat{\mu}^{(\hat{m})}|^2 - 2\langle \varepsilon, \hat{\mu}^{(\hat{m})} - \mu \rangle + \text{Pen}(p_{\hat{m}}) \leq |\mu - \mu^{(m)}|^2 - 2\langle \varepsilon, \mu^{(m)} - \mu \rangle + \text{Pen}(p_m).$$

Umstellen liefert die Behauptung.  $\square$

**Satz 3.19 (Orakelungleichung zur Modellwahl)** *Betrachte das datenerzeugende Modell  $Y = \mu + \varepsilon$ ,  $\mu \in \mathbb{R}^n$ ,  $\varepsilon \sim N(0, \sigma^2 E_n)$  und für lineare Unterräume  $S_m \subseteq \mathbb{R}^n$ ,*

$m = 1, \dots, M$ , mit  $\dim(S_m) = p_m$  die penalisierte Modellwahl  $\hat{m}$ , wobei  $\text{Pen}(p_m) \geq K\sigma^2(p_m + 1)$  für ein  $K > 1$  gelte. Weiter seien  $\hat{\mu}^{(m)} = \Pi_{S_m} Y$ ,  $\mu^{(m)} = \Pi_{S_m} \mu$  und  $\kappa \in (0, \sqrt{K} - 1)$ .

(i) Für jedes  $\tau > 0$  gilt mit der Wahrscheinlichkeit  $1 - e^{-\tau/2} \sum_{m=1}^M e^{-p_m \kappa^2/2}$  die Orakelungleichung:

$$|\hat{\mu}^{(\hat{m})} - \mu|^2 \leq C(K, \kappa) \left( \min_{m=1, \dots, M} (|\mu^{(m)} - \mu|^2 + \text{Pen}(p_m)) + \sigma^2 \tau \right) \quad (3.3)$$

mit einer Konstanten  $C(K, \kappa) > 0$ , die nur von  $K$  und  $\kappa$  abhängt, wobei  $C(K, \kappa) \rightarrow \infty$  für  $\kappa \rightarrow \sqrt{K} - 1$ .

(ii) Es gilt

$$\mathbb{E}[|\hat{\mu}^{(\hat{m})} - \mu|^2] \leq \tilde{C}(K, \kappa) \left( \min_{m=1, \dots, M} (|\mu^{(m)} - \mu|^2 + \text{Pen}(p_m)) + \sigma^2 \sum_{m=1}^M e^{-p_m \kappa^2/2} \right)$$

mit einer Konstanten  $\tilde{C}(K, \kappa) > 0$ , die nur von  $K$  und  $\kappa$  abhängt, wobei  $\tilde{C}(K, \kappa) \rightarrow \infty$  für  $\kappa \rightarrow \sqrt{K} - 1$ .

**Beweis** Der Beweis von (i) erfolgt in vier Schritten:

*Schritt 1: Anwenden der Fundamentalungleichung.* Aus Lemma 3.18 folgt für beliebiges  $m \in \{1, \dots, M\}$

$$|\hat{\mu}^{(\hat{m})} - \mu|^2 \leq |\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) - \text{Pen}(p_{\hat{m}}) + 2\langle \varepsilon, \hat{\mu}^{(\hat{m})} - \mu^{(m)} \rangle.$$

Wegen  $\hat{\mu}^{(\hat{m})} - \mu^{(m)} \in \text{span}(S_{\hat{m}}, \mu^{(m)}) =: S_{\hat{m}}^*$  mit  $\dim S_{\hat{m}}^* \leq p_{\hat{m}} + 1$  folgt

$$\langle \varepsilon, \hat{\mu}^{(\hat{m})} - \mu^{(m)} \rangle \leq |\hat{\mu}^{(\hat{m})} - \mu^{(m)}| \sup_{s \in S_{\hat{m}}^*} \frac{|\langle \varepsilon, s \rangle|}{|s|}.$$

Bezeichnen wir die Orthogonalprojektion auf  $S_{\hat{m}}^*$  mit  $\Pi_{S_{\hat{m}}^*}$ , dann wird das Supremum bei  $s = \Pi_{S_{\hat{m}}^*} \varepsilon$  mit dem maximalen Wert  $|\Pi_{S_{\hat{m}}^*} \varepsilon|$  angenommen. Wir erhalten

$$|\hat{\mu}^{(\hat{m})} - \mu|^2 \leq |\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) - \text{Pen}(p_{\hat{m}}) + 2|\hat{\mu}^{(\hat{m})} - \mu^{(m)}| |\Pi_{S_{\hat{m}}^*} \varepsilon|.$$

*Schritt 2: Konzentrationsungleichung für  $\chi^2(d)$ .* Sei  $Z_d \sim \chi^2(d)$ , das heißt,  $Z_d \stackrel{d}{=} \sum_{i=1}^d X_i^2$  für i.i.d.  $X_i \sim N(0, 1)$ . Wegen  $\mathbb{E}[e^{uX_i^2}] = (1 - 2u)^{-1/2}$  für jedes  $u \in (0, 1/2)$  erhalten wir für  $\rho > 1$

$$\mathbb{P}(Z_d \geq \rho d) \leq \mathbb{E}[e^{uZ_d} e^{-u\rho d}] = (1 - 2u)^{-\frac{d}{2}} e^{-u\rho d} = e^{-\frac{d}{2}(\rho - 1 - \log \rho)},$$

wobei für den letzten Schritt  $u = (\rho - 1)/(2\rho)$  gewählt wurde. Aufgrund von  $\log(1 + t) \leq t$  für  $t \geq 0$  folgt daraus

$$\mathbb{P}(Z_d \geq (1 + \kappa + \sqrt{\tau/d})^2 d)$$

$$\begin{aligned}
&\leq \exp\left(-\frac{d}{2}\left[(\kappa + \sqrt{\tau/d})^2 + 2(\kappa + \sqrt{\tau/d}) - 2\log(1 + \kappa + \sqrt{\tau/d})\right]\right) \\
&\leq \exp\left(-\frac{d}{2}(\kappa + \sqrt{\tau/d})^2\right) \leq \exp\left(-\frac{d}{2}\kappa^2 - \frac{\tau}{2}\right).
\end{aligned}$$

*Schritt 3: Kontrolle des stochastischen Fehlers.* Es gilt  $Z_{d_m} := \sigma^{-2}|\Pi_{S_m^*}\varepsilon|^2 \sim \chi^2(d_m)$  für ein  $d_m \in \{p_m, p_m + 1\}$  und jedes  $m$ . Insbesondere folgt  $\mathbb{E}[|\Pi_{S_m^*}\varepsilon|^2] = \sigma^2 d_m$  für jedes deterministische  $m$ . Wir normieren und korrigieren daher  $|\Pi_{S_m^*}\varepsilon|$  um einen Term dieser Größe. Anschließend schätzen wir  $\widehat{m}$  durch das Maximum in  $\{1, \dots, M\}$  ab, sodass wir die Abhängigkeit zwischen  $\widehat{m}$  und  $\varepsilon$  im Term  $\Pi_{S_{\widehat{m}}^*}\varepsilon$  auflösen:

$$\begin{aligned}
|\Pi_{S_{\widehat{m}}^*}\varepsilon| &\leq \sigma\left((\kappa + 1)\sqrt{p_{\widehat{m}} + 1} + \frac{1}{\sigma}|\Pi_{S_{\widehat{m}}^*}\varepsilon| - (\kappa + 1)\sqrt{p_{\widehat{m}} + 1}\right) \\
&\leq \sigma\left((\kappa + 1)\sqrt{p_{\widehat{m}} + 1} + \max_{m=1, \dots, M} \left(\frac{1}{\sigma}|\Pi_{S_m^*}\varepsilon| - (\kappa + 1)\sqrt{d_m}\right)\right)
\end{aligned}$$

Wir definieren nun das gute Ereignis

$$\mathcal{G} := \left\{ \max_{1 \leq m \leq M} \left( \frac{1}{\sigma}|\Pi_{S_m^*}\varepsilon| - (\kappa + 1)\sqrt{d_m} \right) \leq \sqrt{\tau} \right\}.$$

Aus der Subadditivität von  $\mathbb{P}$  und Schritt 2 angewendet auf  $Z_{d_m} = \sigma^{-2}|\Pi_{S_m^*}\varepsilon|^2$  folgt

$$\begin{aligned}
\mathbb{P}(\mathcal{G}^c) &= \mathbb{P}\left(\max_{1 \leq m \leq M} \left( \frac{1}{\sigma}|\Pi_{S_m^*}\varepsilon| - (\kappa + 1)\sqrt{d_m} \right) \geq \sqrt{\tau}\right) \\
&\leq \sum_{m=1}^M \mathbb{P}\left(\frac{1}{\sigma}|\Pi_{S_m^*}\varepsilon| - (\kappa + 1)\sqrt{d_m} \geq \sqrt{\tau}\right) \\
&\leq \sum_{m=1}^M \mathbb{P}\left(\sqrt{Z_{d_m}} \geq \left(1 + \kappa + \sqrt{\frac{\tau}{d_m}}\right)\sqrt{d_m}\right) \\
&\leq \sum_{m=1}^M \exp\left(-\frac{\kappa^2 d_m}{2} - \frac{\tau}{2}\right) \leq e^{-\tau/2} \sum_{m=1}^M e^{-p_m \kappa^2/2}.
\end{aligned}$$

Damit gilt auf  $\mathcal{G}$  mit der gesuchten Wahrscheinlichkeit

$$|\Pi_{S_{\widehat{m}}^*}\varepsilon| \leq \sigma\left((\kappa + 1)\sqrt{p_{\widehat{m}} + 1} + \sqrt{\tau}\right) \leq \frac{\kappa + 1}{\sqrt{K}}\sqrt{\text{Pen}(p_{\widehat{m}})} + \sigma\sqrt{\tau}.$$

*Schritt 4: Umordnung der Terme.* Auf dem Ereignis  $\mathcal{G}$  erhalten wir durch Kombination von Schritt 1 und Schritt 3:

$$\begin{aligned}
|\widehat{\mu}^{(\widehat{m})} - \mu|^2 &\leq |\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) - \text{Pen}(p_{\widehat{m}}) \\
&\quad + 2|\widehat{\mu}^{(\widehat{m})} - \mu^{(m)}| \left( \frac{\kappa + 1}{\sqrt{K}}\sqrt{\text{Pen}(p_{\widehat{m}})} + \sigma\sqrt{\tau} \right).
\end{aligned}$$

Wir zerlegen nun  $|\widehat{\mu}^{(\widehat{m})} - \mu^{(m)}| \leq |\widehat{\mu}^{(\widehat{m})} - \mu| + |\mu^{(m)} - \mu|$  und verwenden zweimal  $2AB \leq \eta A^2 + \eta^{-1} B^2$ ,  $\eta > 0$ , sodass wir

$$\begin{aligned} |\widehat{\mu}^{(\widehat{m})} - \mu| \left( \frac{\kappa + 1}{\sqrt{K}} \sqrt{\text{Pen}(p_{\widehat{m}})} + \sigma \sqrt{\tau} \right) &\leq (\eta_1^{-1} + \eta_2^{-1}) |\widehat{\mu}^{(\widehat{m})} - \mu|^2 \\ &\quad + \eta_1 \frac{(1 + \kappa)^2}{K} \text{Pen}(p_{\widehat{m}}) + \eta_2 \sigma^2 \tau \end{aligned}$$

für  $\eta_1, \eta_2 > 0$  erhalten. Behandelt man den Summanden mit  $|\mu^{(m)} - \mu|$  analog, dann folgt für  $\eta_3, \eta_4 > 0$

$$\begin{aligned} |\widehat{\mu}^{(\widehat{m})} - \mu|^2 &\leq (1 + \eta_3^{-1} + \eta_4^{-1}) |\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) + (\eta_2 + \eta_4) \tau \sigma^2 \\ &\quad + \left( (\eta_1 + \eta_3) \frac{(1 + \kappa)^2}{K} - 1 \right) \text{Pen}(p_{\widehat{m}}) + (\eta_1^{-1} + \eta_2^{-1}) |\widehat{\mu}^{(\widehat{m})} - \mu|^2. \end{aligned}$$

Aufgrund der Annahme  $\frac{K}{(1+\kappa)^2} > 1$  finden wir nun  $\eta_1, \eta_2, \eta_3 > 0$ , sodass  $\eta_1^{-1} + \eta_2^{-1} < 1$  und  $\eta_1 + \eta_3 = \frac{K}{(1+\kappa)^2}$  erfüllt sind. Außerdem können wir  $\eta_4 = 1$  wählen. Umstellen liefert auf  $\mathcal{G}$ , da  $m$  beliebig war,

$$|\widehat{\mu}^{(\widehat{m})} - \mu|^2 \leq C(\eta_1, \eta_2, \eta_3) \min_{m=1, \dots, M} (|\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) + \sigma^2 \tau).$$

Wir folgern abschließend (ii). Setze  $t^* := C(K, \kappa) \min_{m=1, \dots, M} (|\mu - \mu^{(m)}|^2 + \text{Pen}(p_m))$ . Dann impliziert (i)

$$\begin{aligned} \mathbb{E}[|\widehat{\mu}^{(\widehat{m})} - \mu|^2] &= \int_0^\infty \mathbb{P}(|\widehat{\mu}^{(\widehat{m})} - \mu|^2 > t) dt \\ &\leq \int_0^{t^*} \mathbb{P}(|\widehat{\mu}^{(\widehat{m})} - \mu|^2 > t) dt \\ &\quad + \int_0^\infty \mathbb{P}(|\widehat{\mu}^{(\widehat{m})} - \mu|^2 > t^* + C(K, \kappa) \sigma^2 \tau) C(K, \kappa) \sigma^2 \tau dt \\ &\leq t^* + C(K, \kappa) \sigma^2 \sum_{m=1}^M e^{-p_m \kappa^2 / 2} \int_0^\infty e^{-\tau / 2} d\tau. \end{aligned}$$

Die Behauptung folgt, da das letzte Integral endlich ist.  $\square$

*Bemerkung 3.20 (Interpretation der Orakelungleichung)*

1. Das  $m^*$ , mit dem das Minimum auf der rechten Seite in (3.3) angenommen wird, nennt man auch *Orakelmodell*. Aus der Bias-Varianz-Zerlegung folgt

$$\mathbb{E}[|\widehat{\mu}^{(m)} - \mu|^2] = |\mu^{(m)} - \mu|^2 + \sigma^2 p_m \approx |\mu^{(m)} - \mu|^2 + \text{Pen}(p_m) / K$$

für  $\text{Pen}(d_m) = K \sigma^2 (d_m + 1)$ . Damit liegt  $|\mu^{(m^*)} - \mu|^2 + \text{Pen}(d_{m^*})$  nahe am *Orakelfehler*  $\min_m \mathbb{E}[|\widehat{\mu}^{(m)} - \mu|^2]$ . Mit  $\tau = d_{m^*}$  ist der Restterm  $\sigma^2 \tau$  von der Ordnung  $\text{Pen}(d_{m^*})$  (oder kleiner) und kann damit durch das Minimum in (3.3) abgeschätzt

werden. Für dieses  $\tau$  ist die Wahrscheinlichkeit, mit der die Orakelungleichung gilt, gegeben durch

$$1 - \sum_{m=1}^M e^{-p_m \kappa^2 / 2} e^{-d_{m^*} / 2}.$$

Falls asymptotisch  $n \rightarrow \infty$ ,  $d_{m^*} \rightarrow \infty$  (also  $\mu$  in keinem  $S_m$  exakt liegt) und diese Summe gleichmäßig in  $M$  beschränkt ist, dann konvergiert diese Wahrscheinlichkeit exponentiell schnell in  $d_{m^*}$  gegen eins.

2. In Teil (ii) kann man für festes  $M$  und  $n \rightarrow \infty$  im Fall des BIC mit  $\text{Pen}(p_m) = \log(n)p_m\sigma^2$  auch  $\kappa$  zunehmend größer wählen, sodass der hintere Term sehr klein wird.

Im Satz 3.19 wird nicht gefordert, dass die Modelle geordnet sind. In der multiplen Regression können daher alle  $2^p$  Untermodelle betrachtet werden. Allerdings muss dann  $\kappa$  hinreichend groß gewählt werden, um eine große Wahrscheinlichkeit sicher zu stellen. In der Tat zeigt sich bei dieser „full subset“-Variablenselektion, dass in der Praxis AIC und BIC nicht so gut funktionieren wie größere  $\text{Pen}(p_m)$ -Penalisierungen. Eine weitergehende Analyse von penalisierten Modellwahlverfahren ist in Massart (2007) zu finden.

### 3.3 Aufgaben

- 4.1 Zeigen Sie, dass die Kullback-Leibler-Divergenz nicht symmetrisch und somit keine Metrik ist. Zeigen Sie, dass auch die symmetrisierte Version  $\text{KL}(\mathbb{P}|\mathbb{Q}) + \text{KL}(\mathbb{Q}|\mathbb{P})$  keine Metrik ist.
- 4.2 Bestimmen Sie die Kullback-Leibler-Divergenz  $\text{KL}(\text{N}(\mu_1, \Sigma_1) | \text{N}(\mu_2, \Sigma_2))$  zwischen zwei  $d$ -dimensionalen Normalverteilungen mit  $\mu_1, \mu_2 \in \mathbb{R}^d$  und symmetrischen, positiv-semidefiniten Matrizen  $\Sigma_1, \Sigma_2 \in \mathbb{R}^{d \times d}$ .
- 4.3 Der *Totalvariationsabstand* zweier Wahrscheinlichkeitsmaße  $\mathbb{P}$  und  $\mathbb{Q}$  auf einem messbaren Raum  $(X, \mathcal{F})$  ist definiert als

$$\text{TV}(\mathbb{P}, \mathbb{Q}) := \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)|.$$

- (a) Zeigen Sie

$$\text{TV}(\mathbb{P}, \mathbb{Q}) = \frac{1}{2} \int |p - q| d\nu$$

für ein dominierendes Maß  $\nu$  und die Radon-Nikodym-Dichten  $p = \frac{d\mathbb{P}}{d\nu}$ ,  $q = \frac{d\mathbb{Q}}{d\nu}$ .

- (b) Zeigen Sie  $\text{TV}(\mathbb{P}, \mathbb{Q}) \leq \sqrt{\text{KL}(\mathbb{P}|\mathbb{Q})/2}$ .

*Hinweis:* Betrachten Sie die Funktion  $h(z) := z \log(z) - z + 1$ ,  $z > 0$ , mit stetiger Fortsetzung in Null. Zeigen Sie für alle  $z \geq 0$

$$\left(\frac{4}{3} + \frac{2}{3}z\right)h(z) \geq (z-1)^2.$$

Verwenden Sie anschließend die Cauchy-Schwarz-Ungleichung.

- (c) Prüfen Sie, ob folgende Folgen  $(\mathbb{P}_n)_{n \in \mathbb{N}}$  auf  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  schwach, im Totalvariationsabstand oder in der Kullback-Leibler-Divergenz gegen einen geeigneten Grenzwert  $\mathbb{P}$  konvergieren (betrachten Sie sowohl  $\text{KL}(\mathbb{P}_n|\mathbb{P})$  als auch  $\text{KL}(\mathbb{P}|\mathbb{P}_n)$ ):

- (i)  $\mathbb{P}_n = \delta_{1/n}$ , Dirac-Maß in  $1/n$ ;  
 (ii)  $\mathbb{P}_n = (1 - \frac{1}{n})\nu + \frac{1}{n}\delta_1$ , wobei  $\nu$  das Maß der  $N(0, 1)$ -Verteilung ist.

- 4.4 Wir verwenden polynomielle Untermodelle im nicht parametrischen Regressionsmodell

$$Y_i = f\left(\frac{i}{n}\right) + \varepsilon_i, \quad i = 1, \dots, 100,$$

mit  $\varepsilon_1, \dots, \varepsilon_n \stackrel{i.i.d.}{\sim} N(0, 1)$  und für die (wahren) Funktionen

- (a)  $f: [0, 1] \rightarrow \mathbb{R}, f(x) = \sin(\frac{3}{2}\pi x)$ ,  
 (b)  $f: [0, 1] \rightarrow \mathbb{R}, f(x) = x^2$ .

Verwenden Sie jeweils in 10.000 unabhängigen Simulationen AIC und BIC zur Wahl des Grades  $p-1$  der linearen Modelle  $f(x) = f_k(x) = \beta_0 + \dots + \beta_{p-1}x^{p-1}$ . Bestimmen Sie in allen vier Fällen die Monte-Carlo-Approximation des Vorhersagefehlers  $\mathbb{E}[|Y - X\hat{\beta}|^2]$  und stellen Sie die Wahl von  $p$  jeweils in einem Boxplot dar.

- 4.5 Die reellwertige Zufallsvariable  $X_p$  sei  $\chi^2(p)$ -verteilt mit  $p \in \mathbb{N}$  Freiheitsgraden. Bestimmen Sie eine möglichst scharfe Abschätzung der Wahrscheinlichkeit

$$\mathbb{P}(X_p - \mathbb{E}[X_p] \geq \sqrt{\text{Var}(X_p)\kappa}) \quad \text{für } \kappa > 0.$$

Gehen Sie wie folgt vor:

- (i) Berechnen Sie  $\mathbb{E}[X_p]$ ,  $\text{Var}(X_p)$  und  $\mathbb{E}[\exp(\lambda X_p)]$  für  $\lambda > 0$ .  
 (ii) Verwenden Sie die Markov-Ungleichung, um eine Abschätzung der gewünschten Wahrscheinlichkeit zu erhalten.  
 (iii) Wählen Sie  $\lambda$  optimal.

- 4.6 Das empirische Skalarprodukt sei definiert als  $\langle f, g \rangle_n := \frac{1}{n} \sum_{i=1}^n f(\frac{i}{n})g(\frac{i}{n})$  bzw.  $\langle x, g \rangle_n := \frac{1}{n} \sum_{i=1}^n x_i g(\frac{i}{n})$  für Funktionen  $f, g$  auf  $[0, 1]$  und einen Vektor  $x \in \mathbb{R}^n$ . Die empirische Norm ist gegeben durch  $\|f\|_n^2 := \langle f, f \rangle_n$ . Sei  $(\varphi_k)_{k=1}^n \subseteq L^2([0, 1])$  eine Orthonormalbasis bezüglich  $\langle \cdot, \cdot \rangle_n$ . Für  $f(x) = \sum_{k=1}^n \alpha_k \varphi_k(x), x \in [0, 1]$ , mit Koeffizienten  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$  betrachten wir das Regressionsmodell

$$Y_i = f\left(\frac{i}{n}\right) + \varepsilon_i, \quad i = 1, \dots, n,$$

mit unabhängig und identisch  $N(0, \sigma^2)$ -verteilten Fehlern  $\varepsilon_i, i = 1, \dots, n$ .

- (a) Weisen Sie nach, dass  $\hat{\alpha}_k := \langle Y, \varphi_k \rangle_n$  der Maximum-Likelihood-Schätzer von  $\alpha_k$  für alle  $k = 1, \dots, n$  ist.
- (b) Bestimmen Sie für  $m \in \{1, \dots, n\}$  den Fehler

$$\mathbb{E}[\|f - \hat{f}_m\|_n^2] \quad \text{für} \quad \hat{f}_m(x) := \sum_{k=1}^m \hat{\alpha}_k \varphi_k(x).$$

Auf welches Minimierungsproblem führt eine optimale Wahl von  $m$ ?

- (c) Nehmen Sie an, dass es ein  $s > 0$  und  $0 < c < C$  gibt, sodass  $ck^{-s} \leq \alpha_k \leq Ck^{-s}$  gilt. Wie wächst das optimale  $m$  in  $n$ ?
- 4.7 Betrachten Sie die linearen Modelle  $Y = X^{(p)}\beta^{(p)} + \varepsilon$  mit  $\varepsilon \sim N(0, \sigma^2 E_n)$ ,  $\beta^{(p)} \in \mathbb{R}^p$  und den Designmatrizen  $X^{(p)} \in \mathbb{R}^{n \times p}$ . Leiten Sie die explizite Darstellung des Leave-one-out-Kleinste-Quadrate-Schätzers her.

- 4.8 Betrachten Sie die mittleren Jannuartemperaturen in Berlin-Dahlem in den letzten 300 Jahren (siehe Beispiel 2.52). Verwenden Sie sowohl die polynomiale Regression mit AIC- und BIC-Wahl für den Polynomgrad als auch die Lasso-Methode. Untersuchen Sie die Abhängigkeit vom Tuning-Parameter. Liefern alle Verfahren ähnliche Resultate?

Die beiden folgenden Aufgaben beleuchten das Optimierungsproblem des Lasso-Schätzers genauer. Zu ihrer Lösung sind der Subdifferentialkalkül und die Karush-Kuhn-Tucker-Bedingungen aus der nichtlinearen Optimierung hilfreich.

- 4.9 Betrachten Sie das Regressionsmodell  $Y = \mu + \varepsilon$  mit  $\mu \in \mathbb{R}^n$ , Beobachtungsfehlern  $\varepsilon \in \mathbb{R}^n$  und eine Designmatrix  $X \in \mathbb{R}^{n \times p}$ . Zeigen Sie, dass das Lasso-Optimierungsproblem äquivalent durch

$$\hat{\beta}' = \arg \min_{\beta': |\beta'| \leq s} \frac{1}{n} |Y - X\beta'|^2$$

für einen Sparsity-Parameter  $s > 0$  beschrieben werden kann. Welche Beziehung besteht zwischen  $s$  und dem Penalisierungsparameter  $\lambda$  des Lasso-Schätzers?

- 4.10 Betrachten Sie das lineare Modell  $Y = X\beta + \varepsilon$  mit  $X \in \mathbb{R}^{n \times p}$ ,  $\beta \in \mathbb{R}^p$  und  $\varepsilon \sim N(0, \sigma^2 I_n)$ . Es sei  $G(\beta) := \frac{1}{n} |Y - X\beta|^2$ .

- (a) Kann es mehr als eine Lösung des Optimierungsproblems

$$\min_{\beta \in \mathbb{R}^p} \{G(\beta) + \lambda |\beta|_1\} \quad (*)$$

geben? Begründen Sie Ihre Antwort.

- (b) Zeigen Sie, dass eine notwendige und hinreichende Bedingung an eine Lösung  $\hat{\beta}$  von (\*) gegeben ist durch

$$\nabla_j G(\hat{\beta}) = -\text{sign}(\hat{\beta}_j) \lambda, \quad \text{falls } \hat{\beta}_j \neq 0,$$

$$|\nabla_j G(\widehat{\beta})| \leq \lambda, \quad \text{falls } \widehat{\beta}_j = 0.$$

- (c) Beweisen Sie im Fall, dass (\*) keine eindeutige Lösung besitzt: Gilt  $\nabla G_j(\widehat{\beta}) < \lambda$  für eine Lösung  $\widehat{\beta}$  von (\*), dann gilt  $\widehat{\beta}_j = 0$  für alle Lösungen. Insbesondere hängt die aktive Menge  $\{j : \widehat{\beta}_j \neq 0, j = 1, \dots, p\}$  nicht von der Lösung von (\*) ab.



# Anhang A

## Konzepte der Wahrscheinlichkeitstheorie

### A.1 Grundbegriffe der Maßtheorie und Stochastik

Die Konzepte der Maß- und Wahrscheinlichkeitstheorie bilden die Grundlage der Stochastik und damit insbesondere der Statistik. Der Vollständigkeit halber wollen wir hier die grundlegenden Begriffe und Konzepte zusammenfassen, ohne zu sehr ins Detail zu gehen. Wir verzichten daher in diesem Kapitel auf Beweise.

Für eine grundlegende Einführung sei beispielsweise auf Küchler (2016) verwiesen. Eine tiefgehende und umfangreiche Darstellung der Maßtheorie findet man in Elstrodt (2005). Als Einführung in die Stochastik kann man Georgii (2007) empfehlen. Ein darüber hinausgehendes und weitreichendes Lehrbuch zur Stochastik ist Klenke (2008).

Wir modellieren die Gesamtheit aller möglichen Ausgänge eines Zufallsexperiments durch eine nichtleere Menge  $\Omega$ . Teilmengen von  $\Omega$  beschreiben Ereignisse. Ein Ereignis  $A \subseteq \Omega$  tritt ein, wenn das Ergebnis  $\omega \in \Omega$  des Zufallsexperiments in  $A$  liegt. Die Potenzmenge  $\mathcal{P}(\Omega)$  ist definiert als die Menge aller Teilmengen von  $\Omega$ . Um später jedem Ereignis in konsistenter Art und Weise eine Wahrscheinlichkeit zuordnen zu können, ist im Allgemeinen die gesamte Potenzmenge zu groß. Alle uns interessierenden Ereignisse fassen wir daher in einem Mengensystem  $\mathcal{A}$  aus  $\Omega$  zusammen. Je größer  $\mathcal{A}$  ist, desto genauere Aussagen können wir über den Ausgang des Zufallsexperiments treffen.

**Definition A.1** Für  $\Omega \neq \emptyset$  heißt ein Mengensystem  $\mathcal{A} \subseteq \mathcal{P}(\Omega)$   **$\sigma$ -Algebra über  $\Omega$** , falls folgende Eigenschaften erfüllt sind:

1.  $\Omega \in \mathcal{A}$ ,
2. für alle  $A \in \mathcal{A}$  gilt auch  $A^c \in \mathcal{A}$  und
3. für alle  $A_1, A_2, A_3, \dots \in \mathcal{A}$  gilt auch  $\bigcup_{i \in \mathbb{N}} A_i \in \mathcal{A}$ .

Das Paar  $(\Omega, \mathcal{A})$  heißt **messbarer Raum**.

*Beispiel A.2* Es sei  $\Omega$  eine nichtleere Menge.

- (a)  $\{\emptyset, \Omega\}$  und die Potenzmenge  $\mathcal{P}(\Omega)$  sind  $\sigma$ -Algebren.

(b) Ist  $\mathcal{E} \subseteq \mathcal{P}(\Omega)$  ein System von Teilmengen von  $\Omega$ , dann ist

$$\sigma(\mathcal{E}) := \bigcap \{ \mathcal{A} : \mathcal{A} \text{ ist } \sigma\text{-Algebra aus } \Omega \text{ mit } \mathcal{E} \subseteq \mathcal{A} \}$$

die kleinste  $\sigma$ -Algebra, die  $\mathcal{E}$  umfasst.  $\sigma(\mathcal{E})$  heißt die von  $\mathcal{E}$  **erzeugte  $\sigma$ -Algebra**.

(c) Ist  $\Omega$  mit einer Metrik  $d$  versehen, dann ist die **Borel- $\sigma$ -Algebra**  $\mathcal{B}(\Omega)$  definiert als die kleinste  $\sigma$ -Algebra, die alle offenen Teilmengen von  $\Omega$  enthält:

$$\mathcal{B}(\Omega) = \sigma(\{O \mid O \subseteq \Omega \text{ offen}\})$$

Die Elemente dieser  $\sigma$ -Algebra heißen **Borel-Mengen**.

Die Borel-Mengen sind das (für uns) wichtigste Beispiel einer  $\sigma$ -Algebra. Wenn nichts anderes gefordert oder angenommen wurde, dann verwenden wir für metrische Räume stets die Borel- $\sigma$ -Algebra.

Maße sind Abbildungen von einer  $\sigma$ -Algebra nach  $[0, \infty]$ , die jeder Menge aus der  $\sigma$ -Algebra eine „Größe“ zuordnen, beispielsweise Flächen oder Volumina. Es ist intuitiv, dass diese Abbildungen nichtnegativ und additiv sein sollten.

**Definition A.3** Auf einem messbaren Raum  $(\Omega, \mathcal{A})$  heißt eine nichtnegative Abbildung  $\mu: \mathcal{A} \rightarrow [0, \infty]$  **Maß**, falls

1.  $\mu(\emptyset) = 0$ , die leere Menge also das Maß null hat, und
2.  $\mu$   $\sigma$ -additiv ist, das heißt für paarweise disjunkte Mengen  $A_1, A_2, A_3, \dots \in \mathcal{A}$  mit  $A_i \cap A_j = \emptyset$  für alle  $i \neq j$  gilt

$$\mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \sum_{i \in \mathbb{N}} \mu(A_i).$$

Das Tripel  $(\Omega, \mathcal{A}, \mu)$ , bestehend aus einer nichtleeren Menge  $\Omega$ , einer  $\sigma$ -Algebra  $\mathcal{A}$  und einem Maß  $\mu$  auf  $\mathcal{A}$ , heißt **Maßraum**. Die Mengen  $A \in \mathcal{A}$  mit  $\mu(A) = 0$  heißen  **$\mu$ -Nullmengen**.

Aus der  $\sigma$ -Additivität von Maßen folgt insbesondere die  $\sigma$ -Subadditivität für beliebige Mengen  $A_i \in \mathcal{A}$ ,  $i \in \mathbb{N}$ :

$$\mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) \leq \sum_{i \in \mathbb{N}} \mu(A_i)$$

In der Stochastik werden Maße verwendet, um den Ereignissen Wahrscheinlichkeiten zuzuordnen. Das beschränkt Wertebereich von Wahrscheinlichkeitsmaßen in natürlicher Weise auf  $[0, 1]$ .

**Definition A.4** Auf einem messbaren Raum  $(\Omega, \mathcal{A})$  heißt ein Maß  $\mu$   **$\sigma$ -endlich**, falls eine abzählbare Folge  $(A_n)_{n \in \mathbb{N}}$  von Mengen aus  $\mathcal{A}$  existiert, sodass  $\bigcup_{n \in \mathbb{N}} A_n = \Omega$  sowie  $\mu(A_n) < \infty$  für alle  $n \in \mathbb{N}$  gilt. Ist sogar  $\mu(\Omega) < \infty$ , heißt  $\mu$  **endlich**. Im Spezialfall  $\mu(\Omega) = 1$  wird  $\mu$  **Wahrscheinlichkeitsmaß** und  $(\Omega, \mathcal{A}, \mu)$  **Wahrscheinlichkeitsraum** genannt.

**Beispiel A.5** Wir betrachten einen beliebigen messbaren Raum  $(\Omega, \mathcal{A})$  für eine nicht-leere Menge  $\Omega$ .

1. Das **Zählmaß** auf  $\mathcal{A}$  ist definiert als

$$\mu: \mathcal{A} \rightarrow [0, \infty], \quad \mu(A) = \begin{cases} |A|, & \text{falls } A \text{ endlich ist,} \\ +\infty, & \text{falls } A \text{ unendlich ist.} \end{cases}$$

Jeder messbaren Menge  $A \in \mathcal{A}$  wird also die Anzahl der Elemente in  $A$  zugeordnet.  $\mu$  ist genau dann  $\sigma$ -endlich, wenn  $\Omega$  abzählbar ist. Ist  $\Omega$  endlich, dann ist  $\mathbb{P}(A) := \frac{|A|}{|\Omega|}$ ,  $A \subseteq \Omega$ , ein Wahrscheinlichkeitsmaß, nämlich die Gleichverteilung.

2. Für ein fixiertes Element  $\omega \in \Omega$  ist das **Dirac-Maß** in  $\omega$  definiert via

$$\delta_\omega: \mathcal{A} \rightarrow [0, \infty], \quad \delta_\omega(A) = \begin{cases} 1, & \text{falls } \omega \in A, \\ 0, & \text{sonst.} \end{cases}$$

Man sieht leicht, dass  $\delta_\omega$  ein Wahrscheinlichkeitsmaß ist.

3. Ist  $\Omega = \mathbb{R}$  und  $\mathcal{A} = \mathcal{B}(\mathbb{R})$ , dann ist das (eindeutig bestimmte) Maß  $\lambda$ , für das

$$\lambda((a, b]) = b - a, \quad \text{für alle } a, b \in \mathbb{R}, a < b,$$

gilt, das **Lebesgue-Maß** auf  $\mathbb{R}$ . Analog kann das Lebesgue-Maß auf  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  über Volumina von Quadern eindeutig bestimmt werden.

Meist betrachten wir Wahrscheinlichkeitsmaße aus einer der beiden folgenden Klassen.

**Definition A.6** Es sei  $(\Omega, \mathcal{A}, \mathbb{P})$  ein Wahrscheinlichkeitsraum.

1. Ist  $\Omega$  eine endliche oder abzählbar unendliche Menge und  $\mathcal{A} = \mathcal{P}(\Omega)$ , so nennen wir  $(\Omega, \mathcal{A}, \mathbb{P})$  einen **diskreten Wahrscheinlichkeitsraum**. In diesem Fall existiert stets eine **Zähldichte** von  $\mathbb{P}$  gegeben durch  $p: \Omega \rightarrow [0, 1]$  mit  $\sum_{\omega \in \Omega} p(\omega) = 1$  und  $p(\omega) = \mathbb{P}(\{\omega\})$ . Es gilt

$$\mathbb{P}(A) = \sum_{\omega \in A} p(\omega) \quad \text{für alle } A \in \mathcal{A}.$$

2. Ist  $\Omega = \mathbb{R}$  versehen mit der  $\sigma$ -Algebra  $\mathcal{A} = \mathcal{B}(\mathbb{R})$  und existiert eine Funktion  $f: \mathbb{R} \rightarrow [0, \infty)$  mit  $\int_{\mathbb{R}} f(x) dx = 1$  und

$$\mathbb{P}((a, b]) = \int_a^b f(x) dx \quad \text{für alle } a, b \in \mathbb{R}, a < b,$$

so heißt  $f$  **Wahrscheinlichkeitsdichte** von  $\mathbb{P}$  oder kurz **Dichte** von  $\mathbb{P}$ . Wir sprechen in diesem Fall von einer **stetigen Wahrscheinlichkeitsverteilung**.

Man beachte, dass durch  $\mathbb{P}((a, b]) = \int_a^b f(x) dx$  das Maß  $\mathbb{P}$  bereits eindeutig durch seine Dichte festgelegt wird (was aus dem Eindeutigkeitssatz aus der Maßtheorie

folgt). Zudem ist der Fall von stetigen Verteilungen leicht auf den mehrdimensionalen Fall  $\mathbb{R}^d$  zu übertragen. Wichtige Beispiele für diskrete und stetige Verteilungen werden in den Abschnitten A.2 und A.3 beschrieben.

**Definition A.7** In einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  gilt ein Ereignis  $A \in \mathcal{A}$   **$\mathbb{P}$ -fast sicher** (kurz:  $\mathbb{P}$ -f.s.), wenn  $\mathbb{P}(A) = 1$  gilt. Die Eigenschaft  $E(\omega)$  sei für die Elemente  $\omega \in \Omega$  sinnvoll. Dann sagt man, die Eigenschaft  $E$  gilt  **$\mathbb{P}$ -fast überall** auf  $\Omega$  (kurz:  $\mathbb{P}$ -f.ü.), wenn es eine  $\mathbb{P}$ -Nullmenge  $N \in \mathcal{A}$  gibt, sodass alle  $\omega \in N^c$  die Eigenschaft  $E$  haben. Man sagt auch, dass  $E$  für  **$\mathbb{P}$ -fast alle** ( $\mathbb{P}$ -f.a.)  $\omega \in \Omega$  gilt.

Nachdem wir einen Wahrscheinlichkeitsraum definiert haben, wollen wir als Nächstes Zufallsvariablen einführen.

**Definition A.8** Seien  $(\Omega, \mathcal{A})$  und  $(\mathcal{X}, \mathcal{F})$  zwei messbare Räume. Eine Abbildung  $X: \Omega \rightarrow \mathcal{X}$  heißt  $(\mathcal{A}, \mathcal{F})$ -**messbar**, falls

$$\forall A \in \mathcal{F}: X^{-1}(A) := \{\omega \in \Omega \mid X(\omega) \in A\} \in \mathcal{A}.$$

Hierbei bezeichnet  $X^{-1}(A)$  das Urbild von  $A$  unter  $X$ . Gehen die  $\sigma$ -Algebren  $\mathcal{A}$  und  $\mathcal{F}$  aus dem Kontext eindeutig hervor, nennen wir  $X$  kurz **messbar**. Messbare Abbildungen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{F}, \mathbb{P})$  werden **Zufallsvariablen** genannt.

Ist eine Zufallsvariable  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ -wertig, so sprechen wir auch kurz von einer reellen Zufallsvariable.

*Beispiel A.9* Gegeben sei der messbare Raum  $(\Omega, \mathcal{A})$ .

1. Gilt  $n \in \mathbb{N}$ ,  $A_1, \dots, A_n \in \mathcal{A}$  und  $a_1, \dots, a_n \in \mathbb{R} \setminus \{0\}$ , so ist

$$X(\omega) = \sum_{k=1}^n a_k \mathbb{1}_{A_k}(\omega), \quad \omega \in \Omega,$$

$\mathcal{A}$ -messbar. Funktionen dieser Gestalt nennen wir **einfache Funktionen**. Man kann zeigen, dass jede messbare nichtnegative Funktion durch eine monotone Folge einfacher Funktionen approximiert werden kann.

2. Jede stetige Funktion  $X: \mathbb{R} \rightarrow \mathbb{R}$  ist Borel-messbar, das heißt  $(\mathcal{B}(\mathbb{R}), \mathcal{B}(\mathbb{R}))$ -messbar. Die umgekehrte Implikation gilt nicht, da die Dirichlet-Funktion  $X(y) = \mathbb{1}_{\mathbb{Q}}(y)$ ,  $y \in \mathbb{R}$ , messbar, aber nirgendwo stetig ist.
3. Es seien  $\Omega$  eine Menge,  $(\mathcal{X}, \mathcal{F})$  ein messbarer Raum und  $X: \Omega \rightarrow \mathcal{X}$  eine Abbildung. Dann ist  $\sigma(X) := X^{-1}(\mathcal{F}) = \{X^{-1}(A) : A \in \mathcal{F}\}$  die kleinste  $\sigma$ -Algebra  $\mathcal{A}$  auf  $\Omega$ , sodass  $X$  eine  $(\mathcal{A}, \mathcal{F})$ -messbare Abbildung ist. Wir nennen  $\sigma(X)$  die von  $X$  **erzeugte  $\sigma$ -Algebra**.

Die folgende Eigenschaft erlaubt es Stochastikern, sich auf Wahrscheinlichkeitsmaße auf dem Ergebnisraum eines Zufallsexperiments zu beschränken, statt den gesamten zugrunde liegenden Wirkmechanismus beschreiben zu müssen.

**Lemma A.10** Es seien  $(\Omega, \mathcal{A}, \mu)$  ein Maßraum,  $(X, \mathcal{F})$  ein messbarer Raum und  $X: \Omega \rightarrow X$  eine  $(\mathcal{A}, \mathcal{F})$ -messbare Abbildung. Durch

$$\mu^X(A) := \mu(X^{-1}(A)), \quad A \in \mathcal{F},$$

ist auf  $\mathcal{F}$  ein Maß  $\mu^X$  definiert, das als von  $X$  **induziertes Maß** oder als **Bildmaß** von  $X$  bezeichnet wird. Ist  $\mu$  ein Wahrscheinlichkeitsmaß, so ist auch  $\mu^X$  ein Wahrscheinlichkeitsmaß.

**Definition A.11** Für eine Zufallsvariable  $X$  auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  wird das induzierte Maß  $\mathbb{P}^X$  auch **Verteilung der Zufallsvariable**  $X$  genannt, und wir schreiben dann  $X \sim \mathbb{P}^X$ .

Analog zu Definition A.6 sprechen wir von diskret bzw. stetig verteilten Zufallsvariablen  $X$ , falls ihre Bildmaße diskret mit Zähldichte  $p^X$  bzw. stetig mit Dichte  $f^X$  sind.

*Bemerkung A.12* Wir betrachten einen Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$ , einen messbaren Raum  $(X, \mathcal{F})$  und eine Zufallsvariable  $X: (\Omega, \mathcal{A}) \rightarrow (X, \mathcal{F})$ . Für  $A \in \mathcal{F}$  sind folgende Schreibweisen in der Stochastik üblich und zweckmäßig:

$$\begin{aligned} \{X \in A\} &:= \{\omega \in \Omega \mid X(\omega) \in A\} = X^{-1}(A) \quad \text{sowie} \\ \mathbb{P}(X \in A) &:= \mathbb{P}(\{X \in A\}) = \mathbb{P}(\{\omega \in \Omega \mid X(\omega) \in A\}) = \mathbb{P}^X(A). \end{aligned}$$

Die Verteilung einer reellen Zufallsvariable kann (eindeutig) durch ihre Verteilungsfunktion beschrieben werden. Letztere ist wie folgt definiert.

**Definition A.13** Sei  $X$  eine reelle Zufallsvariable auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$ . Die **Verteilungsfunktion** von  $X$  ist definiert als

$$F^X: \mathbb{R} \rightarrow [0, 1], \quad F^X(x) := \mathbb{P}(X \leq x) = \mathbb{P}^X((-\infty, x]).$$

Für eine diskret auf einer abzählbaren Teilmenge  $S$  von  $\mathbb{R}$  verteilten Zufallsvariable ist ihre Verteilungsfunktion eine stückweise konstante Treppenfunktion mit Sprüngen auf  $S$ . Für stetig verteilte Zufallsvariablen mit Wahrscheinlichkeitsdichte  $f^X$  ist die Verteilungsfunktion stetig und es gilt  $F^X(x) = \int_{-\infty}^x f^X(y) dy$ .

Verteilungsfunktionen sind stets monoton wachsend, rechtsstetig, das heißt  $F^X(x) = \lim_{y \downarrow x} F^X(y)$  für alle  $x \in \mathbb{R}$ , und es existieren die linken Grenzwerte  $F^X(x-) := \lim_{y \uparrow x} F^X(y)$ . Zudem gilt  $\lim_{x \rightarrow -\infty} F^X(x) = 0$  sowie  $\lim_{x \rightarrow \infty} F^X(x) = 1$ . Diese Eigenschaften implizieren insbesondere, dass eine (verallgemeinerte) Inverse von  $F^X$  existiert.

**Definition A.14** Sei  $F: \mathbb{R} \rightarrow \mathbb{R}$  eine monoton wachsende Funktion, wobei wir den Definitionsbereich um  $\pm\infty$  erweitern:

$$F(-\infty) := \lim_{x \rightarrow -\infty} F(x), \quad F(\infty) = \lim_{x \rightarrow \infty} F(x)$$

Die **verallgemeinerte Inverse**  $F^-: \mathbb{R} \rightarrow \overline{\mathbb{R}} = [-\infty, \infty]$  von  $F$  ist definiert als

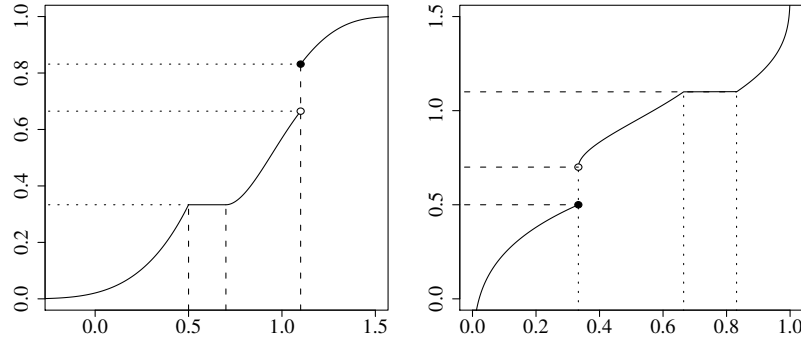


Abb. A.1 Eine Verteilungsfunktion (links) und die zugehörige Quantilfunktion (rechts)

$$F^-(y) := \inf\{x \in \mathbb{R} : F(x) \geq y\}, \quad y \in \mathbb{R},$$

mit der Konvention  $\inf \emptyset := \infty$ . Falls  $F: \mathbb{R} \rightarrow [0, 1]$  eine Verteilungsfunktion ist, heißt  $F^-: [0, 1] \rightarrow \mathbb{R}$  **Quantilfunktion** von  $F$ .

**Definition A.15** Es sei  $X$  eine Zufallsvariable auf dem Wahrscheinlichkeitsraum  $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbb{P})$  mit zugehöriger Verteilungsfunktion  $F^X$ . Für jedes  $p \in (0, 1)$  ist die Menge aller  $p$ -**Quantile** von  $F^X$  gegeben durch

$$\{x \in \mathbb{R} : \mathbb{P}(X \leq x) \geq p \text{ und } \mathbb{P}(X \geq x) \geq 1 - p\}.$$

Die Menge der  $p$ -Quantile ist ein abgeschlossenes Intervall mit der Obergrenze

$$\sup\{x \in \mathbb{R} : \mathbb{P}(X \geq x) \geq 1 - p\}$$

und der Untergrenze  $(F^X)^-(p)$ . Falls  $F^X$  invertierbar ist, fallen Ober- und Untergrenze zusammen und das  $p$ -Quantil ist eindeutig. Im Verlauf des Buches wird klar, dass Quantile für die Konstruktion und Kalibrierung vieler statistischer Methoden von grundlegender Bedeutung sind.

Wir kommen nun zu zwei zentralen Begriffen der Wahrscheinlichkeitstheorie, nämlich bedingten Wahrscheinlichkeiten und Unabhängigkeit. Erstere formalisieren die Intuition, dass sich Wahrscheinlichkeiten von Ereignissen ändern, wenn über den Ausgang eines Zufallsexperiments bereits eine Teilinformation vorhanden ist.

**Definition A.16** Es seien  $(\Omega, \mathcal{A}, \mathbb{P})$  ein Wahrscheinlichkeitsraum und  $A, B \in \mathcal{A}$  Ereignisse mit  $\mathbb{P}(B) > 0$ . Dann ist die **bedingte Wahrscheinlichkeit** von  $A$  gegeben  $B$  definiert als

$$\mathbb{P}(A|B) := \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

Es ist leicht nachzuprüfen, dass  $\mathbb{P}(\cdot|B)$  tatsächlich ein Wahrscheinlichkeitsmaß ist, falls  $\mathbb{P}(B) > 0$ . Falls  $\mathbb{P}(A|B) = \mathbb{P}(A)$  für zwei Ereignisse  $A, B \in \mathcal{A}$  gilt, so enthält

$B$  keinerlei Information über  $A$ . Die Ereignisse  $A$  und  $B$  sind also unabhängig. Umstellen zeigt, dass  $\mathbb{P}(A|B) = \mathbb{P}(A)$  äquivalent zu  $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$  ist.

**Definition A.17** Auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  heißen zwei Ereignisse  $A, B \in \mathcal{A}$  **(stochastisch) unabhängig**, falls  $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$  gilt. Eine Familie von Ereignissen  $(A_i)_{i \in I}$  mit nichtleerer Indexmenge  $I \neq \emptyset$  heißt **(stochastisch) unabhängig**, falls  $\mathbb{P}(\bigcap_{i \in J} A_i) = \prod_{i \in J} \mathbb{P}(A_i)$  für jede endliche Teilmenge  $J \subseteq I$ .

Der folgende Satz enthält zwei elementare, aber wichtige Eigenschaften von bedingten Wahrscheinlichkeiten.

**Satz A.18 (Gesetz der totalen Wahrscheinlichkeit, Bayes-Formel)** Es sei  $(\Omega, \mathcal{A}, \mathbb{P})$  ein Wahrscheinlichkeitsraum und  $(B_i)_{i \in I}$  eine höchstens abzählbare Folge paarweiser disjunkter Mengen aus  $\mathcal{A}$  mit  $\bigcup_{i \in I} B_i = \Omega$  und  $\mathbb{P}(B_i) > 0$  für alle  $i \in I$ .

(i) Für jedes Ereignis  $A \in \mathcal{A}$  gilt das Gesetz der totalen Wahrscheinlichkeit

$$\mathbb{P}(A) = \sum_{i \in I} \mathbb{P}(B_i) \mathbb{P}(A|B_i).$$

(ii) Für jedes  $A \in \mathcal{A}$  mit  $\mathbb{P}(A) > 0$  und alle  $i \in I$  gilt die Bayes-Formel

$$\mathbb{P}(B_i|A) = \frac{\mathbb{P}(B_i) \mathbb{P}(A|B_i)}{\sum_{k \in I} \mathbb{P}(B_k) \mathbb{P}(A|B_k)} = \frac{\mathbb{P}(B_i) \mathbb{P}(A|B_i)}{\mathbb{P}(A)}.$$

**Definition A.19** Eine Familie  $(X_i)_{i \in I}$  von  $(S_i, \mathcal{S}_i)$ -wertigen Zufallsvariablen heißt **unabhängig**, falls für jede beliebige Wahl  $A_i \in \mathcal{S}_i$  die Familie von Ereignissen  $(X_i \in A_i)_{i \in I}$  unabhängig ist.

Ein Vektor oder eine Folge  $X = (X_i)_{i \in I}$  von  $S_i$ -wertigen Zufallsvariablen  $X_i$ ,  $i \in I$ , nimmt Werte im Produktraum  $\prod_{i \in I} S_i$  an. Wollen wir die Verteilung  $X$  auf diesem beschreiben, müssen wir den Produktraum zunächst mit einer  $\sigma$ -Algebra versehen.

**Definition A.20** Es seien  $(\Omega_i, \mathcal{A}_i, \mathbb{P}_i)$ ,  $i \in I$ , Wahrscheinlichkeitsräume für eine beliebige nichtleere Indexmenge  $I \neq \emptyset$ . Die **Produkt- $\sigma$ -Algebra** über dem kartesischen Produkt  $\Omega := \prod_{i \in I} \Omega_i$  ist definiert als

$$\mathcal{A} := \bigotimes_{i \in I} \mathcal{A}_i := \sigma(\{\pi_i^{-1}(A_i) | i \in I, A_i \in \mathcal{A}_i\}),$$

wobei  $\pi_i: \Omega \rightarrow \Omega_i$  die  $i$ -te Koordinatenprojektion bezeichnet. Gilt für ein Wahrscheinlichkeitsmaß  $\mathbb{P}$  auf  $\mathcal{A}$

$$\mathbb{P}\left(\bigcap_{i \in J} \pi_i^{-1}(A_i)\right) = \prod_{i \in J} \mathbb{P}_i(A_i) \quad \text{für alle } J \subseteq I \text{ endlich, } A_i \in \mathcal{A}_i,$$

so heißt  $\mathbb{P}$  **Produktmaß**. Wir schreiben  $\mathbb{P} = \bigotimes_{i \in I} \mathbb{P}_i$ .

Sind  $X_i: \Omega \rightarrow \mathcal{X}_i$  für  $i = 1, \dots, n$  Zufallsvariablen auf einem gemeinsamen Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  mit Werten in  $(\mathcal{X}_i, \mathcal{F}_i)$ , dann ist auch der Vektor  $(X_1, \dots, X_n)^\top$  eine messbare Abbildung, das heißt eine Zufallsvariable:

$$(X_1, \dots, X_n)^\top: (\Omega, \mathcal{A}) \rightarrow \left( \prod_{i=1}^n \mathcal{X}_i, \bigotimes_{i=1}^n \mathcal{F}_i \right)$$

Die Verteilung  $\mathbb{P}^{(X_1, \dots, X_n)}$  von  $(X_1, \dots, X_n)^\top$  wird auch die **gemeinsame Verteilung** der Zufallsvariablen  $X_1, \dots, X_n$  genannt.

**Satz A.21** *Es seien  $X_i: \Omega \rightarrow \mathcal{X}_i$  Zufallsvariablen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  mit Werten in den messbaren Räumen  $(\mathcal{X}_i, \mathcal{F}_i)$  für  $i \in I \neq \emptyset$ . Dann ist die Familie  $(X_i)_{i \in I}$  genau dann unabhängig, wenn für alle endlichen Teilmengen  $J \subseteq I$  die gemeinsame Verteilung von  $(X_i)_{i \in J}$  auf  $(\prod_{i \in J} \mathcal{X}_i, \bigotimes_{i \in J} \mathcal{F}_i)$  das Produktmaß der Marginalverteilungen ist, das heißt wenn  $\mathbb{P}^{(X_i)_{i \in J}} = \bigotimes_{i \in J} \mathbb{P}^{X_i}$ .*

Man kann zeigen, dass für gegebene  $(\Omega_i, \mathcal{A}_i, \mathbb{P}_i)_{i \in I}$  mit beliebiger Indexmenge  $I \neq \emptyset$ , stets ein Produktmaß  $\mathbb{P} = \bigotimes_{i \in I} \mathbb{P}_i$  auf dem Produktraum  $(\prod_{i \in I} \Omega_i, \bigotimes_{i \in I} \mathcal{A}_i)$  existiert. Insbesondere findet man also einen Wahrscheinlichkeitsraum, auf dem beliebig viele unabhängige Zufallsvariablen  $X_i$  mit vorgegeben Randverteilungen  $\mathbb{P}^{X_i} = \mathbb{P}_i$  existieren. Sind  $(X_i)$  unabhängig und identisch verteilt, das heißt  $\mathbb{P}^{X_i} = \mathbb{P}^{X_j}$  für alle  $i, j \in I$  beschreiben wir dies häufig mit **i.i.d.**, was das englische *independent and identically distributed* abkürzt.

Um Momente, wie den Erwartungswert oder die Varianz, von reellwertigen Zufallsvariablen zu definieren, benötigen wir einen geeigneten Integralbegriff auf dem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$ . Die Integrationstheorie von Lebesgue erlaubt eine Integration bezüglich allgemeiner (Wahrscheinlichkeits-)Maße. Die Integral-konstruktion beruht auf drei Schritten. Für einfache Zufallsvariablen

$$X = \sum_{i=1}^n \alpha_i \mathbb{1}_{A_i} \quad \text{mit} \quad n \in \mathbb{N}, \alpha_1, \dots, \alpha_n \in \mathbb{R}, A_1, \dots, A_n \in \mathcal{A},$$

siehe Beispiel A.9, ist der Erwartungswert bzw. das Integral von  $X$  bezüglich  $\mathbb{P}$  definiert als

$$\mathbb{E}[X] := \int_{\Omega} X d\mathbb{P} := \sum_{i=1}^n \alpha_i \mathbb{P}(A_i).$$

Dieser Erwartungswert hängt nicht von der Darstellung von  $X$  ab, ist linear in  $X$  und monoton. Für jede nichtnegative Zufallsvariable  $X$  finden wir im zweiten Schritt stets eine Folge von einfachen Zufallsvariablen  $X_n$ , sodass  $X_n(\omega) \uparrow X(\omega)$  für alle  $\omega \in \Omega$  gilt, siehe Beispiel A.9. Dann definieren wir

$$\mathbb{E}[X] := \lim_{n \rightarrow \infty} \mathbb{E}[X_n].$$

Aufgrund der Monotonie des Erwartungswerts und der Folge  $(X_n)$  ist dieser Grenzwert wohldefiniert und man kann zeigen, dass diese Definition nicht von der Wahl der approximierenden Folge abhängt.

**Definition A.22** Auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  ist die Menge aller endlich integrierbaren Zufallsvariablen gegeben durch

$$\mathcal{L}^1 := \mathcal{L}^1(\Omega, \mathcal{A}, \mathbb{P}) := \{X: \Omega \rightarrow \mathbb{R} \text{ messbar} \mid \mathbb{E}[|X|] < \infty\}.$$

Für  $X \in \mathcal{L}^1$  definieren wir mit  $X_+ := \max(X, 0)$  und  $X_- := \max(-X, 0)$  den **Erwartungswert** als

$$\mathbb{E}[X] := \mathbb{E}[X_+] - \mathbb{E}[X_-] \in \mathbb{R}.$$

Wir schreiben  $\mathbb{E}[X] = \int_{\Omega} X \, d\mathbb{P} = \int_{\Omega} X(\omega) \mathbb{P}(d\omega)$  und  $\int_A X \, d\mathbb{P} = \int_{\Omega} X(\omega) \mathbb{1}_A(\omega) \mathbb{P}(d\omega)$  für  $A \in \mathcal{A}$ .

Analog wird das Integral bezüglich allgemeiner Maße definiert. Folgender Satz fasst die elementarsten Eigenschaften des Erwartungswerts zusammen.

**Satz A.23** Für  $X \in \mathcal{L}^1(\Omega, \mathcal{A}, \mathbb{P})$  gilt:

- (i)  $\mathbb{E}[X] = \int_{\mathbb{R}} x \mathbb{P}^X(dx)$ , insbesondere hängt der Erwartungswert nur von der Verteilung  $\mathbb{P}^X$  von  $X$  ab.
- (ii) Der Erwartungswert ist linear und monoton: Ist  $Y \in \mathcal{L}^1$  eine weitere Zufallsvariable und sind  $\alpha, \beta \in \mathbb{R}$ , so gilt  $\mathbb{E}[\alpha X + \beta Y] = \alpha \mathbb{E}[X] + \beta \mathbb{E}[Y]$ . Aus  $X \leq Y$  folgt  $\mathbb{E}[X] \leq \mathbb{E}[Y]$ .
- (iii) Falls  $X, Y \in \mathcal{L}^1$  unabhängig sind, so gilt  $X \cdot Y \in \mathcal{L}^1$  und  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ .

In den Spezialfällen diskret bzw. stetig verteilter Zufallsvariablen ergibt sich:

**Korollar A.24** Es sei  $X$  eine Zufallsvariable auf dem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$ .

- (i) Besitzt  $X$  einen abzählbaren Wertebereich  $X(\Omega) \subseteq \mathbb{R}$ , so gilt  $X \in \mathcal{L}^1$  genau dann, wenn  $\sum_{x \in X(\Omega)} |x| \mathbb{P}(X = x)$  endlich ist. In diesem Fall gilt für den Erwartungswert

$$\mathbb{E}[X] = \sum_{x \in X(\Omega)} x \mathbb{P}(X = x).$$

- (ii) Ist  $X$  eine Zufallsvariable mit Dichte  $f^X: \mathbb{R} \rightarrow [0, \infty)$ , so gilt  $X \in \mathcal{L}^1$  genau dann, wenn  $\int_{\mathbb{R}} |x| f^X(x) dx$  endlich ist. In diesem Fall gilt für den Erwartungswert

$$\mathbb{E}[X] = \int_{\mathbb{R}} x f^X(x) dx.$$

Die Darstellung des Erwartungswerts in Abhängigkeit von der Verteilung der Zufallsvariable gilt auch allgemeiner. Ist  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor und  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  Borel-messbar, dann ist die Zufallsvariable  $h(X)$  genau dann in  $\mathcal{L}^1$ , wenn  $\int_{\mathbb{R}^d} |h(x)| \mathbb{P}^X(dx) < \infty$ . In diesem Fall ist der Erwartungswert gegeben durch

$$\mathbb{E}[h(X)] = \int_{\mathbb{R}^d} h(x) \mathbb{P}^X(dx).$$

**Bemerkung A.25** Für eine  $\mathbb{R}^d$ -wertige Zufallsvariable  $X$  auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  ist die Familie der Funktionen  $h_u(x) = \exp(i\langle x, u \rangle)$  für  $u \in \mathbb{R}^d$  und mit der imaginären Einheit  $i = \sqrt{-1}$  ein wichtiger Spezialfall. Der Erwartungswert von  $h_u(X)$  wird durch

$$\begin{aligned} \varphi(u) &:= \mathbb{E}[e^{i\langle u, X \rangle}] := \mathbb{E}[\operatorname{Re} e^{i\langle u, X \rangle}] + i\mathbb{E}[\operatorname{Im} e^{i\langle u, X \rangle}] \\ &= \mathbb{E}[\cos(\langle u, X \rangle)] + i\mathbb{E}[\sin(\langle u, X \rangle)] \end{aligned}$$

definiert. Da Kosinus und Sinus durch eins beschränkt sind, sind letztere Erwartungswerte stets wohldefiniert. Die Funktion  $u \mapsto \varphi(u)$  heißt **charakteristische Funktion**, und man kann beweisen, dass die Verteilung von  $X$  eindeutig durch  $\varphi$  bestimmt wird.

Die Monome  $h(x) = x^p$  führen uns auf folgende Definition:

**Definition A.26** Eine Zufallsvariable  $X$  liegt in  $\mathcal{L}^p$  für  $p > 0$ , falls  $|X|^p \in \mathcal{L}^1$ , also falls  $\mathbb{E}[|X|^p] < \infty$  gilt. In diesem Fall heißt  $\mathbb{E}[|X|^p]$  das  **$p$ -te absolute Moment** von  $X$ . Für  $X \in \mathcal{L}^p$  und  $p \in \mathbb{N}$  heißt  $\mathbb{E}[X^p]$  das  **$p$ -te Moment** von  $X$ .

Man beachte, dass  $\mathcal{L}^q \subseteq \mathcal{L}^p$  für  $0 < p \leq q$  gilt, wobei hierfür essentiell ist, dass  $\mathbb{P}$  ein endliches Maß ist. Von herausragender Bedeutung ist das zentrierte zweite Moment, also die Varianz, da sie die Streuung einer Zufallsvariable um ihren Erwartungswert angibt. Die lineare Abhängigkeit zwischen zwei Zufallsvariablen wird über ihre Korrelation quantifiziert.

**Definition A.27** Für eine Zufallsvariable  $X \in \mathcal{L}^2$  bezeichnet

$$\operatorname{Var}(X) := \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

die **Varianz** von  $X$ . Ihre Wurzel  $\sqrt{\operatorname{Var}(X)}$  heißt **Standardabweichung** von  $X$ . Für  $X, Y \in \mathcal{L}^2$  definiert

$$\operatorname{Cov}(X, Y) := \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

die **Kovarianz** zwischen  $X$  und  $Y$ . Gilt  $\operatorname{Var}(X), \operatorname{Var}(Y) > 0$ , ist die **Korrelation** gegeben durch

$$\rho(X, Y) = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X) \operatorname{Var}(Y)}}.$$

Im Fall  $\operatorname{Cov}(X, Y) = 0$  heißen  $X$  und  $Y$  **unkorreliert**.

**Bemerkung A.28** Diese Definitionen lassen sich leicht auf  $d$ -dimensionale Zufallsvariablen übertragen  $X = (X_1, \dots, X_d)^\top$ : Den Erwartungswert verstehen wir komponentenweise, das heißt

$$\mathbb{E}[X] := (\mathbb{E}[X_1], \dots, \mathbb{E}[X_d])^\top.$$

Die Varianz des Zufallsvektors  $X$  definieren wir als erwarteten euklidischen Abstand vom Erwartungswert:

$$\text{Var}(X) := \mathbb{E}[|X - \mathbb{E}[X]|^2]$$

Die Kovarianzen zwischen allen Komponenten von  $X$  werden in der **Kovarianzmatrix** zusammengefasst:

$$\text{Cov}(X) := \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top] = (\text{Cov}(X_i, X_j))_{i,j=1,\dots,d} \in \mathbb{R}^{d \times d}$$

**Satz A.29** Für reellwertige Zufallsvariablen  $X, Y$  auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  gelten folgende Ungleichungen:

1. Für  $X \in \mathcal{L}^1(\mathbb{P})$  gilt die **Markov-Ungleichung**

$$\mathbb{P}(|X| > \kappa) \leq \frac{\mathbb{E}[|X|]}{\kappa} \quad \text{für alle } \kappa > 0.$$

Gilt  $X \in \mathcal{L}^2(\mathbb{P})$ , folgt die **Tschebyscheff-Ungleichung**

$$\mathbb{P}(|X - \mathbb{E}[X]| > \kappa) \leq \frac{\text{Var}(X)}{\kappa^2} \quad \text{für alle } \kappa > 0.$$

2. Für  $X, Y \in \mathcal{L}^2(\mathbb{R})$  ist  $XY \in \mathcal{L}^1(\mathbb{P})$ , und es gilt die **Cauchy-Schwarz-Ungleichung**

$$\mathbb{E}[|XY|] \leq \mathbb{E}[X^2]^{1/2} \mathbb{E}[Y^2]^{1/2}.$$

Sind allgemeiner  $p, q \geq 1$  derart, dass  $1 = \frac{1}{p} + \frac{1}{q}$  und  $X \in \mathcal{L}^p(\mathbb{P}), Y \in \mathcal{L}^q(\mathbb{P})$ , dann folgt  $XY \in \mathcal{L}^1(\mathbb{P})$ , und es gilt die **Hölder-Ungleichung**

$$\mathbb{E}[|XY|] \leq \mathbb{E}[|X|^p]^{1/p} \mathbb{E}[|Y|^q]^{1/q}.$$

3. Ist  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  konvex und  $\varphi(X) \in \mathcal{L}^1(\mathbb{P})$ , dann gilt die **Jensen-Ungleichung**

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)].$$

Die Momente einer Verteilung beschreiben Eigenschaften wie den Mittelwert, die Streuung, die Schiefe (zentriertes drittes Moment, englisch: *skewness*) oder die Wölbung (zentriertes viertes Moment, englisch: *kurtosis*). Darüber hinaus kann man sich fragen, ob die Folge der Momente  $m_n = \mathbb{E}[X^n]$ ,  $n \in \mathbb{N}$ , für  $X \sim \mathbb{P}$  ein Wahrscheinlichkeitsmaß  $\mathbb{P}$  eindeutig bestimmt, falls alle Momente wohldefiniert und endlich sind. Diese Frage ist als das *Momentenproblem* bekannt und auch von statistischer Relevanz, siehe Momentenmethode in Abschnitt 1.2. Basierend auf der charakteristischen Funktion aus Bemerkung A.25 geben wir exemplarisch folgende Charakterisierung an:

**Satz A.30** Seien  $\mathbb{P}$  ein Wahrscheinlichkeitsmaß auf dem messbaren Raum  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  und  $X \sim \mathbb{P}$  mit existierenden Momenten  $m_n = \mathbb{E}[X^n]$ ,  $n \in \mathbb{N}$ , und charakteristischer Funktion  $\varphi(u) = \mathbb{E}[e^{iuX}]$ ,  $u \in \mathbb{R}$ . Dann sind folgende Eigenschaften äquivalent:

- (i)  $\varphi$  ist analytisch auf einer Umgebung um die Null.
- (ii)  $\varphi$  ist analytisch auf  $\mathbb{R}$ .
- (iii)  $\limsup_{n \rightarrow \infty} (|m_n|/n!)^{1/n} < \infty$ .

Jede der drei Eigenschaften impliziert, dass die Momentenfolge  $(m_n)_{n \in \mathbb{N}}$  das Maß  $\mathbb{P}$  eindeutig bestimmt.

Aus Eigenschaft (iii) erhalten wir das folgende Korollar.

**Korollar A.31** Sei  $\mathbb{P}$  ein Wahrscheinlichkeitsmaß auf  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  mit kompaktem Träger, das heißt, die kleinste abgeschlossene Menge mit Maß eins ist kompakt. Dann ist  $\mathbb{P}$  eindeutig durch seine Momente bestimmt.

Um den Einfluss verschiedener Parameterwahlen auf die Wahrscheinlichkeitsverteilung zu beschreiben, gibt es in statistischen Modellen nicht nur ein Wahrscheinlichkeitsmaß auf dem zugrunde liegenden Raum  $(\Omega, \mathcal{A})$ , sondern mehrere. Um diese zueinander in Beziehung zu setzen, spielt der Satz von Radon-Nikodym eine zentrale Rolle.

**Definition A.32** Es seien  $\mu$  und  $\nu$  Maße auf einem messbaren Raum  $(\Omega, \mathcal{A})$ .  $\nu$  heißt **absolutstetig** bezüglich  $\mu$ , falls jede  $\mu$ -Nullmenge eine  $\nu$ -Nullmenge ist, das heißt, für jedes  $A \in \mathcal{A}$  mit  $\mu(A) = 0$  gilt auch  $\nu(A) = 0$ . Wir schreiben  $\nu \ll \mu$ .

**Satz A.33 (Radon-Nikodym)** Es sei  $(X, \mathcal{A})$  ein messbarer Raum mit einem  $\sigma$ -endlichen Maß  $\mu$  und einem Maß  $\nu$ , das absolutstetig bezüglich  $\mu$  ist. Dann existiert eine messbare Funktion  $f: X \rightarrow [0, \infty]$ , sodass

$$\nu(A) = \int_A f d\mu \quad \text{für alle} \quad A \in \mathcal{A}.$$

$f$  ist  $\mu$ -fast überall eindeutig bestimmt und heißt **Radon-Nikodym-Dichte** von  $\nu$  bezüglich  $\mu$  oder  **$\mu$ -Dichte von  $\nu$** . Wir schreiben  $\frac{d\nu}{d\mu} := f$ .

Mit Blick auf Definition A.6 stellen wir fest, dass uns der Satz von Radon-Nikodym einen einheitlichen Rahmen für diskrete und stetige Verteilungen liefert: Die Zähldichte einer diskreten Verteilung ist gerade die Radon-Nikodym-Dichte bezüglich des Zählmaßes, während die Wahrscheinlichkeitsdichte einer stetigen Verteilung als Radon-Nikodym-Dichte bezüglich des Lebesgue-Maßes aufgefasst werden kann.

Bedingte Wahrscheinlichkeiten  $\mathbb{P}(A | B)$  sind nur für Ereignisse  $B$  mit  $\mathbb{P}(B) > 0$  wohldefiniert. Mithilfe von Radon-Nikodym-Dichten lässt sich dies verallgemeinern.

**Definition A.34** Es seien  $(X, \mathcal{F}, \mu)$  und  $(Y, \mathcal{G}, \nu)$  Maßräume und  $X, Y$  Zufallsvariablen auf  $(\Omega, \mathcal{A}, \mathbb{P})$  mit Werten in  $X$  bzw.  $Y$ , deren gemeinsame Verteilung  $\mathbb{P}^{X,Y}$  auf  $(X \times Y, \mathcal{F} \otimes \mathcal{G})$  eine Dichte  $f^{X,Y}: X \times Y \rightarrow \mathbb{R}$  bezüglich dem Produktmaß  $\mu \otimes \nu$  besitzt. Dann heißt

$$f^{Y|X=x}(y) := \frac{f^{X,Y}(x, y)}{f^X(x)} \quad \text{mit} \quad f^X(x) := \int_Y f(x, z) \nu(dz)$$

für alle  $x \in \mathcal{X}$  mit positiver Randdichte  $f^X(x) > 0$  **bedingte Dichte** von  $Y$  gegeben  $X = x$ . Für alle  $x \in \mathcal{X}$  mit  $f^X(x) = 0$  wird  $f^{Y|X=x}$  beliebig gesetzt, zum Beispiel null. Für messbare Mengen  $A \in \mathcal{G}$  und messbare Funktionen  $\varphi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  bezeichnet

$$\mathbb{P}(Y \in A \mid X = x) := \int_A f^{Y|X=x}(y) \nu(dy)$$

die bedingte Wahrscheinlichkeit für  $Y \in A$ , gegeben  $X = x$ , und

$$\mathbb{E}[\varphi(X, Y) \mid X = x] := \int_{\mathcal{Y}} \varphi(x, y) f^{Y|X=x}(y) \nu(dy)$$

den **bedingten Erwartungswert** von  $\varphi(X, Y)$ , gegeben  $X = x$ , sofern das Integral wohldefiniert ist.

In Analogie zu  $\mathbb{P}(A \cap B) = \mathbb{P}(A \mid B)\mathbb{P}(B)$  gilt also  $f^{X,Y}(x, y) = f^{Y|X=x}(y)f^X(x)$ . Die Formel von der totalen Wahrscheinlichkeit  $\mathbb{P}(A) = \sum_i \mathbb{P}(A \mid B_i)\mathbb{P}(B_i)$  für eine Partition  $\bigcup_i B_i = \Omega$  findet ihre Entsprechung in

$$\mathbb{P}(Y \in A) = \int_{\mathcal{X}} \mathbb{P}(Y \in A \mid X = x) f^X(x) \mu(dx).$$

Allgemeiner gilt auch

$$\mathbb{E}[\varphi(X, Y)] = \int_{\mathcal{X}} \mathbb{E}[\varphi(X, Y) \mid X = x] f^X(x) \mu(dx).$$

Die Bayes-Formel aus Satz A.18 verallgemeinert sich für bedingte Dichten zu folgendem Zusammenhang:

**Satz A.35 (Bayes-Formel)** *In der Situation von Definition A.34 gilt*

$$f^{Y|X=x}(y) = \frac{f^{X|Y=y}(x)f^Y(y)}{\int_{\mathcal{X}} f^{X|Y=z}(x)f^Y(z)dz} = \frac{f^{X|Y=y}(x)f^Y(y)}{f^X(x)} \quad \text{für } \mathbb{P}^X\text{-f.a. } x \in \mathcal{X}.$$

Mit dem Satz von Radon-Nikodym lässt sich eine abstrakte Definition von  $\mathbb{E}[\varphi(X, Y) \mid X = x]$  für beliebige Zufallsvariablen  $X, Y$  mit  $\varphi(X, Y) \in \mathcal{L}^1$  geben, die im Fall einer gemeinsamen Dichte  $f^{X,Y}$  der obigen entspricht. Mit Indikatorfunktionen  $\varphi$  ergeben sich dann auch entsprechende bedingte Wahrscheinlichkeiten. Für die Zwecke dieses Buches reicht obige Definition aus.

Wir schließen dieses Grundlagenkapitel mit den zentralen Sätzen aus der Wahrscheinlichkeitstheorie ab: dem Gesetz der großen Zahlen und dem zentralen Grenzwertsatz. Im Gegensatz zur klassischen Analysis unterscheiden wir in der Stochastik zwischen verschiedenen Konvergenzarten.

**Definition A.36** Es sei  $(\mathcal{X}, d)$  ein metrischer Raum. Eine Folge von Wahrscheinlichkeitsmaßen  $(\mathbb{P}_n)_{n \in \mathbb{N}}$  auf  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$  **konvergiert schwach** gegen ein Wahrscheinlichkeitsmaß  $\mathbb{P}$  auf  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ , wenn für alle stetigen und beschränkten Funktion

$\varphi: \mathcal{X} \rightarrow \mathbb{R}$

$$\lim_{n \rightarrow \infty} \int_{\mathcal{X}} \varphi d\mathbb{P}_n = \int_{\mathcal{X}} \varphi d\mathbb{P}$$

gilt. Wir schreiben  $\mathbb{P}_n \xrightarrow{w} \mathbb{P}$ , wobei das  $w$  für *weak* steht, oder  $\mathbb{P}_n \xrightarrow{d} \mathbb{P}$  mit  $d$  für *distribution*. Eine Folge von Zufallsvariablen  $(X_n)_{n \in \mathbb{N}}$  mit Werten in  $(\mathcal{X}, d)$  **konvergiert in Verteilung** bzw. **schwach** gegen eine Zufallsvariable  $X$  in  $\mathcal{X}$ , falls  $\mathbb{P}^{X_n} \xrightarrow{w} \mathbb{P}^X$ . Wir schreiben kurz  $X_n \xrightarrow{w} X$  oder  $X_n \xrightarrow{d} X$ .

Die schwache Konvergenz der Folge  $(X_n)_{n \in \mathbb{N}}$  gegen die Zufallsvariable  $X$  ist also äquivalent zu

$$\lim_{n \rightarrow \infty} \mathbb{E}[\varphi(X_n)] = \mathbb{E}[\varphi(X)]$$

für alle stetigen und beschränkten Funktionen  $\varphi: \mathcal{X} \rightarrow \mathbb{R}$ . Ist die Folge  $(X_n)_{n \in \mathbb{N}}$  reellwertig, so konvergiert sie genau dann in Verteilung gegen die Zufallsvariable  $X$ , wenn die Folge ihrer Verteilungsfunktionen  $(F^{X_n})_{n \in \mathbb{N}}$  an jeder Stetigkeitsstelle der Verteilungsfunktion  $F^X$  punktweise gegen  $F^X$  konvergiert, das heißt, für alle Stetigkeitsstellen  $x$  von  $F^X$  gilt

$$\lim_{n \rightarrow \infty} F^{X_n}(x) = F^X(x).$$

**Definition A.37** Eine Folge  $(X_n)_{n \in \mathbb{N}}$  von Zufallsvariablen auf einem gemeinsamen Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, \mathbb{P})$  mit Werten in einem normierten Raum  $(\mathcal{X}, |\cdot|)$  konvergiert **stochastisch** bzw. **in Wahrscheinlichkeit** gegen die Zufallsvariable  $X$ , wenn für alle  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0$$

gilt. Wir schreiben kurz  $X_n \xrightarrow{\mathbb{P}} X$ . Die Folge  $(X_n)_{n \in \mathbb{N}}$  konvergiert **fast sicher** gegen die Zufallsvariable  $X$ , wenn

$$\mathbb{P}(\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\}) = 1$$

ist. In diesem Fall schreiben wir  $X_n \xrightarrow{\text{f.s.}} X$ .

Man beachte, dass fast sichere Konvergenz stochastische Konvergenz impliziert und dass aus stochastischer Konvergenz schwache Konvergenz folgt.

**Satz A.38 (Starkes Gesetz der großen Zahlen)** Es sei  $(X_i)_{i \in \mathbb{N}}$  eine Folge von Zufallsvariablen in  $\mathcal{L}^1$  mit demselben Erwartungswert  $\mu = \mathbb{E}[X_i]$ . Weiter sei eine der beiden folgenden Bedingungen erfüllt:

- (i)  $(X_i)_{i \geq 1}$  sind identisch verteilt und paarweise unabhängig.
- (ii)  $(X_i)_{i \geq 1}$  liegen in  $\mathcal{L}^2$ , sind paarweise unkorreliert, und es gilt  $\sup_{i \in \mathbb{N}} \text{Var}(X_i) < \infty$ .

Dann gilt

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{f.s.}} \mu.$$

**Satz A.39 (Zentraler Grenzwertsatz)** Es sei  $(X_i)_{i \geq 1}$  eine Folge unabhängiger und identisch verteilter Zufallsvariablen in  $\mathcal{L}^2(\Omega, \mathcal{A}, \mathbb{P})$  mit  $\mu = \mathbb{E}[X_1]$  und  $\sigma^2 = \text{Var}(X_1)$ . Dann erfüllt ihre standardisierte Summe

$$Z_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} \xrightarrow{d} N(0, 1),$$

wobei  $N(0, 1)$  die Standardnormalverteilung bezeichnet. Insbesondere gilt für  $a < b$  also  $\lim_{n \rightarrow \infty} \mathbb{P}(a < Z_n \leq b) = \Phi(b) - \Phi(a)$  mit der Verteilungsfunktion  $\Phi$  der Standardnormalverteilung.

Ein weiterer, wichtiger Satz ist das *continuous mapping*-Theorem, mit dem man Konsistenz und Grenzwertverteilungseigenschaften übertragen kann.

**Satz A.40 (Continuous Mapping)** Seien die Abbildung  $f: \mathbb{R}^d \mapsto \mathbb{R}^k$  stetig und  $(X_n)_{n \in \mathbb{N}}$  eine Folge von  $d$ -dimensionalen Zufallsvariablen, die schwach, fast sicher bzw. stochastisch gegen  $X \in \mathbb{R}^d$  konvergiert. Dann konvergiert auch  $f(X_n)$  schwach, fast sicher bzw. stochastisch gegen  $f(X)$ .

Folgender Kalkül der stochastischen Ordnung verallgemeinert Landaus  $O$ -Symbol der Analysis und ist oft hilfreich.

**Definition A.41** Für eine Folge reellwertiger Zufallsvariablen  $X_n$  und positive Zahlen  $a_n$  schreiben wir  $X_n = O_{\mathbb{P}}(a_n)$ , falls

$$\lim_{R \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(|X_n| > Ra_n) = 0,$$

und sagen, dass  $X_n$  die **stochastische Ordnung**  $a_n$  besitzt. Im Fall  $X_n = O_{\mathbb{P}}(1)$  nennen wir die Folge  $(X_n)$  **stochastisch beschränkt**.

Ist  $(X_n)$  deterministisch, so gilt  $X_n = O_{\mathbb{P}}(a_n)$  genau dann, wenn  $X_n \leq Ca_n$  für eine Konstante  $C > 0$  gilt, das heißt  $X_n = O(a_n)$ . Für eine Nullfolge  $(a_n)$  folgt aus  $X_n = O_{\mathbb{P}}(a_n)$ , dass  $\mathbb{P}(|X_n| > \varepsilon)$  für jedes  $\varepsilon > 0$  gegen null konvergiert, also  $X_n \xrightarrow{\mathbb{P}} 0$ . Stochastische Beschränktheit einer Folge  $(X_n)$  ist auch als Straffheit von Verteilungen in der Wahrscheinlichkeitstheorie bekannt. Weitere wichtige Eigenschaften werden in dem folgenden Satz zusammengefasst.

**Satz A.42 (Eigenschaften des  $O_{\mathbb{P}}$ -Kalküls)** Für reellwertige Zufallsvariablen  $X_n, Y_n$  und positive Zahlen  $a_n, b_n$  gilt:

1. Aus  $\mathbb{E}[|X_n|^p]^{1/p} \leq Ca_n$  für Konstanten  $C > 0$  und  $p > 0$  folgt  $X_n = O_{\mathbb{P}}(a_n)$  (aus  $L^p$ -Beschränktheit folgt stochastische Beschränktheit).
2. Aus  $a_n^{-1} X_n \xrightarrow{d} Y$  für eine Zufallsvariable  $Y$  folgt  $X_n = O_{\mathbb{P}}(a_n)$  (Konvergenz in Verteilung impliziert stochastische Beschränktheit).
3. Aus  $X_n = O_{\mathbb{P}}(a_n)$ ,  $Y_n = O_{\mathbb{P}}(b_n)$  folgt  $X_n + Y_n = O_{\mathbb{P}}(a_n + b_n)$  und  $X_n Y_n = O_{\mathbb{P}}(a_n b_n)$ , symbolisch:

$$O_{\mathbb{P}}(a_n) + O_{\mathbb{P}}(b_n) = O_{\mathbb{P}}(a_n + b_n), \quad O_{\mathbb{P}}(a_n) O_{\mathbb{P}}(b_n) = O_{\mathbb{P}}(a_n b_n).$$

Wir verwenden die letzte Eigenschaft auch für zufällige Matrizen  $M_n$  und Vektoren  $v_n$  in der Form, dass  $\|M_n\| = O_{\mathbb{P}}(a_n)$  (mit Spektralnorm  $\|\cdot\|$ ) und  $|v_n| = O_{\mathbb{P}}(b_n)$  implizieren  $|M_n v_n| = O_{\mathbb{P}}(a_n b_n)$ , was wegen  $|M_n v_n| \leq \|M_n\| |v_n|$  offensichtlich ist.

## A.2 Diskrete Verteilungen

Im Folgenden sollen häufig auftretende diskrete Verteilungen eingeführt werden. Wir beginnen mit dem einfachsten Fall, bei dem es nur zwei mögliche Versuchsausgänge gibt.

**Definition A.43** Eine Zufallsvariable  $X \in \{0, 1\}$  heißt **Bernoulli-verteilt**, falls  $\mathbb{P}(X = 1) = p$  und  $\mathbb{P}(X = 0) = 1 - p$  für eine Erfolgswahrscheinlichkeit  $p \in [0, 1]$  gilt. Man schreibt  $X \sim \text{Ber}(p)$ .

Für  $X \sim \text{Ber}(p)$  gilt  $\mathbb{E}[X] = p$  und  $\text{Var}(X) = p(1 - p)$ .

Auf endlichen Grundräumen ist die Gleichverteilung die wohl wichtigste, insbesondere mit Blick auf Urnenmodelle.

**Definition A.44** Ist  $\Omega \neq \emptyset$  ein endlicher Grundraum mit Kardinalität  $|\Omega| \in \mathbb{N}$ , dann ist die diskrete **Gleichverteilung**  $U(\Omega)$  gegeben durch die Zähldichte

$$U(\Omega)(\{\omega\}) = \frac{1}{|\Omega|} \quad \text{für alle } \omega \in \Omega.$$

Zur Beschreibung der Anzahl der Erfolge in einer Serie von  $n \in \mathbb{N}$  gleichartigen und unabhängigen Bernoulli-Versuchen nutzt man die Binomialverteilung.

**Definition A.45** Für  $n \in \mathbb{N}$  und  $p \in [0, 1]$  ist die **Binomialverteilung**  $\text{Bin}(n, p)$  durch die Zähldichte

$$\text{Bin}(n, p)(\{k\}) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k \in \{0, 1, \dots, n\},$$

gegeben.

Der Erwartungswert einer binomialverteilten Zufallsvariable  $X \sim \text{Bin}(n, p)$  ist  $np$ . Die Varianz beträgt  $np(1 - p)$ .

Wir verallgemeinern ein Binomialesperiment auf mehrere mögliche Versuchsausgänge. Betrachten wir also wieder  $n \in \mathbb{N}$  unabhängige und gleich verteilten Durchläufe, wobei jeweils  $s \in \mathbb{N}$  verschiedene Versuchsausgänge  $\{1, \dots, s\}$  möglich sind und  $j \in \{1, \dots, s\}$  mit Wahrscheinlichkeit  $p_j$  eintritt. Bezeichne  $X_j$  die Anzahl der Durchgänge aus  $n$  Versuchen, bei denen  $j$  aufgetreten ist, dann ist der Vektor  $X := (X_1, \dots, X_s)$  multinomialverteilt.

**Definition A.46** Für  $n, s \in \mathbb{N}$  und  $p_1, \dots, p_s \in [0, 1]$  mit  $\sum_{j=1}^s p_j = 1$  ist eine Zufallsvariable  $X := (X_1, \dots, X_s) \sim \text{Mult}(n, p_1, \dots, p_s)$  **multinomial-verteilt**, falls für  $k = (k_1, \dots, k_s) \in \mathcal{X}$  mit

$$\mathcal{X} := \{k = (k_1, \dots, k_s) : k_1, \dots, k_s \in \mathbb{N}_0 \text{ mit } k_1 + \dots + k_s = n\}$$

gilt:

$$\mathbb{P}(X = k) = \frac{n!}{k_1! \cdot k_2! \cdot \dots \cdot k_s!} \cdot p_1^{k_1} \cdot p_2^{k_2} \cdot \dots \cdot p_s^{k_s}$$

Wir nennen  $\frac{n!}{k_1! \cdot k_2! \cdot \dots \cdot k_s!}$  **Multinomialkoeffizient**.

Das Binomial- und das Multinomialmodell können wir zur Beschreibung von Urnenmodellen nutzen, wobei die Erfolgswahrscheinlichkeiten dem relativen Anteil der Kugeln in einer bestimmten Farbe entsprechen. Da diese Wahrscheinlichkeiten in jedem Versuchsdurchgang, also in jeder Ziehung, gleich sind, handelt es sich um eine Ziehung mit Zurücklegen. Was passiert ohne Zurücklegen? In einer Grundgesamtheit gäbe es  $N \in \mathbb{N}$  Elemente, wobei jedes Element nur eine von zwei möglichen Ausprägungen hat (zum Beispiel Erfolg/Misserfolg oder rot/blau). Sind  $M \leq N$  Elemente mit der gewünschten Eigenschaft (Erfolg oder rot) in der Grundgesamtheit, so gibt die hypergeometrische Verteilung die Wahrscheinlichkeit an, beim Ziehen ohne Zurücklegen in der Stichprobe  $k \leq M$  Elemente mit der gewünschten Eigenschaft zu finden.

**Definition A.47** Für  $N, M, n \in \mathbb{N}$  mit  $M \leq N$  besitzt eine Zufallsvariable  $X \sim \text{Hyp}(N, M, n)$  die **hypergeometrische Verteilung**, falls

$$\mathbb{P}(X = k) = \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}}, \quad k \in \{0, 1, \dots, n\}.$$

Der Erwartungswert einer  $\text{Hyp}(N, M, n)$ -verteilten Zufallsvariable  $X$  ist  $\mathbb{E}[X] = n \frac{M}{N}$ .

*Bemerkung A.48* Der Unterschied zwischen der hypergeometrischen und der Binomialverteilung ist das Zurücklegen. Dieser Effekt ist bei einem relativ kleinen Stichprobenumfang  $n$  im Vergleich zu  $M$ , das heißt wenn  $n/M \rightarrow 0$ , vernachlässigbar gering. Tatsächlich gilt für alle  $0 \leq k \leq n$

$$\begin{aligned} \binom{M}{k} &= \frac{M^k}{k!} \frac{M(M-1) \cdots (M-k+1)}{M^k} \\ &= \frac{M^k}{k!} \left(1 - \frac{(M-1)}{M}\right) \cdots \left(1 - \frac{(M-k+1)}{M}\right) = \frac{M^k}{k!} (1 + o(1)) \quad \text{für } \frac{n}{M} \rightarrow 0. \end{aligned}$$

Für  $X \sim \text{Hyp}(N, M, n)$  gilt also punktweise für jedes  $k \in \{0, 1, \dots, n\}$ :

$$\begin{aligned} \mathbb{P}(X = k) &= \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}} = \binom{n}{k} \frac{M^k (N-M)^{n-k}}{N^n} (1 + o(1)) \\ &= \binom{n}{k} \left(\frac{M}{N}\right)^k \left(1 - \frac{M}{N}\right)^{n-k} (1 + o(1)). \end{aligned}$$

Folglich kann man für  $n/M \rightarrow 0$  die hypergeometrische Verteilung  $\text{Hyp}(N, M, n)$  durch die einfacher zu handhabende Binomialverteilung  $\text{Bin}(n, \frac{M}{N})$  approximieren.

Abschließend führen wir noch eine Verteilung ein, die auch als Verteilung seltener Ereignisse bezeichnet wird (wegen des Grenzwertsatzes von Poisson).

**Definition A.49** Die **Poisson-Verteilung**  $\text{Poiss}(\lambda)$  mit dem Intensitätsparameter  $\lambda > 0$  ist durch die Zähldichte

$$\text{Poiss}(\lambda)(\{k\}) = \frac{\lambda^k}{k!} \exp(-\lambda), \quad k \in \mathbb{N}_0,$$

gegeben.

Für eine Zufallsvariable  $X \sim \text{Poiss}(\lambda)$  gilt  $\mathbb{E}[X] = \text{Var}(X) = \lambda$ .

### A.3 Stetige Verteilungen

In diesem Kapitel wollen wir einige wichtige stetige Wahrscheinlichkeitsverteilungen zusammenfassend einführen. Diese sind insbesondere durch ihre Dichte charakterisiert. Wir betrachten dabei stets die Borel- $\sigma$ -Algebra auf dem jeweiligen (Teil-)Raum der reellen Zahlen. Im einfachsten Fall ist die Dichte konstant.

**Definition A.50** Sei  $\Omega = [a, b]$  für  $a < b \in \mathbb{R}$ . Mit  $U([a, b])$  bezeichnet man die **Gleichverteilung** auf  $[a, b]$ . Die Dichte  $f$  einer Zufallsvariable  $X \sim U([a, b])$  ist für alle  $x \in \mathbb{R}$  gegeben durch

$$f(x) = \frac{1}{b-a} \mathbb{1}_{[a,b]}(x).$$

Der Erwartungswert und die Varianz einer  $U([a, b])$ -verteilten Zufallsvariable  $X$  sind  $\mathbb{E}[X] = \frac{a+b}{2}$  und  $\text{Var}(X) = \frac{1}{12}(b-a)^2$ .

Im zentralen Grenzwertsatz ist uns bereits die Normalverteilung begegnet, die von fundamentaler Bedeutung ist.

**Definition A.51** Eine Zufallsvariable  $X$  auf  $\mathbb{R}$  mit der Wahrscheinlichkeitsdichte  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right),$$

heißt **normalverteilt** (oder auch **Gauß-verteilt**) mit den Parametern  $\mu \in \mathbb{R}$  und  $\sigma^2 > 0$ . Wir schreiben  $X \sim N(\mu, \sigma^2)$ . Die Verteilungsfunktion der Standardnormalverteilung  $N(0, 1)$  wird mit

$$\Phi(x) := \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt, \quad x \in \mathbb{R},$$

bezeichnet.

Für eine Zufallsvariable  $X \sim N(\mu, \sigma)$  geben die beiden Parameter gerade den Erwartungswert  $\mathbb{E}[X] = \mu$  und die Varianz  $\text{Var}(X) = \sigma^2$  an. Die Normalverteilung konzentriert sich sehr stark um ihren Mittelwert: Für  $X \sim N(0, 1)$  gilt nämlich

$$\mathbb{P}(|X| \geq t) = 2 \int_t^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \leq \frac{2}{\sqrt{2\pi}} \int_t^\infty \frac{x}{t} e^{-x^2/2} dx = \frac{\sqrt{2}}{\sqrt{\pi}t} e^{-t^2/2}. \quad (\text{A.1})$$

**Definition A.52** Für  $d \in \mathbb{N}$  ist ein  $d$ -dimensionaler reeller Zufallsvektor  $X = (X_1, \dots, X_d)$  **mehrdimensional** (oder auch **multivariat**) **normalverteilt** mit Erwartungsvektor  $\mu \in \mathbb{R}^d$  und positiv definiter Kovarianzmatrix  $\Sigma \in \mathbb{R}^{d \times d}$ , wenn sie eine Dichtefunktion der Form

$$f_X(x) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right), \quad x \in \mathbb{R}^d,$$

besitzt. Man schreibt  $X \sim N(\mu, \Sigma)$ .

Für einen Zufallsvektor  $X \sim N(\mu, \Sigma)$  erhalten wir  $\mathbb{E}[X_i] = \mu_i$  und  $\text{Cov}(X_i, X_j) = \Sigma_{i,j}$  für alle  $i, j \in \{1, \dots, d\}$  bzw.  $\mathbb{E}[X] = \mu$  und  $\text{Cov}(X) = \Sigma$  in Vektor- bzw. Matrixnotation.

*Beispiel A.53* Für einen Zufallsvektor  $X \sim N(\mu, \Sigma)$  mit  $\mu \in \mathbb{R}^d$  und positiv definiter Kovarianzmatrix  $\Sigma \in \mathbb{R}^{d \times d}$  gilt für die Lineartransformation  $AX + b \in \mathbb{R}^p$  mit  $A \in \mathbb{R}^{p \times d}$ ,  $b \in \mathbb{R}^p$ , dass

$$AX + b \sim N(A\mu + b, A\Sigma A^\top).$$

Sind  $X_1, \dots, X_m \sim N(0, 1)$  unabhängig, dann kann man auch die Dichte von  $X := \sum_{i=1}^m X_i^2$  explizit bestimmen und erhält die folgende Verteilung:

**Definition A.54** Eine reellwertige Zufallsvariable  $X$  ist  $\chi^2$ -**verteilt** mit  $m$  Freiheitsgraden, geschrieben  $X \sim \chi^2(m)$ , wenn ihre Verteilung durch die Dichte

$$f_X(x) = \frac{1}{2^{m/2} \Gamma(m/2)} x^{m/2-1} e^{-x/2} \mathbb{1}_{\{x>0\}}$$

gegeben ist, wobei  $\Gamma(x) := \int_0^\infty t^{x-1} \exp(-t) dt$  die Gammafunktion bezeichnet.

Es folgen weitere stetige Verteilungen.

**Definition A.55** Mit  $\text{Exp}(\lambda)$ ,  $\lambda > 0$ , bezeichnet man die **Exponentialverteilung**. Die Dichte  $f$  einer Zufallsvariable  $X \sim \text{Exp}(\lambda)$  ist für alle  $x \in \mathbb{R}$  gegeben durch

$$f(x) = \lambda e^{-\lambda x} \mathbb{1}_{\{x \geq 0\}}.$$

Der Erwartungswert und die Varianz der  $\text{Exp}(\lambda)$ -Verteilung sind  $\mathbb{E}[X] = \frac{1}{\lambda}$  und  $\text{Var}(X) = \frac{1}{\lambda^2}$ . Eine Verallgemeinerung von Exponential- und  $\chi^2$ -Verteilung ist durch die Gammaverteilung gegeben.

**Definition A.56** Eine Zufallsvariable  $X$  mit der Wahrscheinlichkeitsdichte  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \mathbb{1}_{\{x>0\}},$$

heißt **gammaverteilt** mit den positiven, reellen Parametern  $\alpha, \beta > 0$ , geschrieben  $X \sim \Gamma(\alpha, \beta)$ .

Wir erhalten  $\text{Exp}(\lambda) = \Gamma(1, \lambda)$  und  $\chi^2(m) = \Gamma(m/2, 1/2)$ . Für eine gammaverteilte Zufallsvariable  $X$  gilt  $\mathbb{E}[X] = \alpha/\beta$  und  $\text{Var}(X) = \alpha/\beta^2$ .

**Definition A.57** Die **Betaverteilung**  $\text{Beta}(\alpha, \beta)$  ist eine stetige Wahrscheinlichkeitsverteilung auf  $([0, 1], \mathcal{B}([0, 1]))$  mit den Parametern  $\alpha, \beta > 0$ . Ihre Wahrscheinlichkeitsdichte ist

$$f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \mathbb{1}_{\{x \in [0,1]\}}$$

mit der **Betafunktion**

$$B(\alpha, \beta) := \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx.$$

Der Erwartungswert und die Varianz einer Zufallsvariable  $X \sim \text{Beta}(\alpha, \beta)$  sind

$$\mathbb{E}[X] = \frac{\alpha}{\alpha+\beta} \quad \text{und} \quad \text{Var}(X) = \frac{\alpha\beta}{(\alpha+\beta+1)(\alpha+\beta)^2}.$$

*Beispiel A.58* Sortieren wir  $n \in \mathbb{N}$  unabhängige und  $U([0, 1])$ -verteilte Zufallsvariablen  $U_1, \dots, U_n$  der Größe nach, erhalten wir die Ordnungsstatistiken  $U_{(1)} \leq \dots \leq U_{(n)}$ . Für die  $j$ -te Ordnungsstatistik gilt  $U_{(j)} \sim \text{Beta}(j, n-j+1)$  und insbesondere  $\mathbb{E}[U_{(j)}] = \frac{j}{n+1}$ . Der Abstand zweier Ordnungsstatistiken ist  $U_{(j)} - U_{(j-1)}$  ist  $\text{Beta}(1, n)$ -verteilt.

# Literaturverzeichnis

- Abbott, B. P., Abbott, R., Abbott, T. D., Abernathy, M. R., Acernese, F., Ackley, K., et al. (2016). Binary black hole mergers in the first advanced ligo observing run. *Phys. Rev. X*, 6:041015.
- Dickhaus, T. (2014). *Simultaneous Statistical Inference with Applications in the Life Sciences*. Springer-Verlag, Berlin Heidelberg.
- Elstrodt, J. (2005). *Maß- und Integrationstheorie*. Springer-Lehrbuch. Grundwissen Mathematik. Springer-Verlag, Berlin, 4. edition.
- Fahrmeir, L., Kneib, T., and Lang, S. (2009). *Regression. Modelle, Methoden und Anwendungen*. Springer-Verlag Berlin Heidelberg.
- Georgii, H.-O. (2007). *Stochastik. Einführung in die Wahrscheinlichkeitstheorie und Statistik*. de Gruyter Lehrbuch. Walter de Gruyter & Co., Berlin, 3. edition.
- Klenke, A. (2008). *Wahrscheinlichkeitstheorie*. Springer-Verlag Berlin Heidelberg, 2. edition.
- Küchler, U. (2016). *Maßtheorie für Statistiker. Grundlagen der Stochastik*. Springer-Verlag Berlin Heidelberg.
- Lehmann, E. L. and Casella, G. (1998). *Theory of point estimation*. Springer Texts in Statistics. Springer-Verlag, New York, 2. edition.
- Lehmann, E. L. and Romano, J. P. (2005). *Testing statistical hypotheses*. Springer Texts in Statistics. Springer, New York, third edition.
- Massart, P. (2007). *Concentration inequalities and model selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard.
- Vershynin, R. (2018). *High-dimensional probability. An introduction with applications in data science*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge. With a foreword by Sara van de Geer.
- Wainwright, M. J. (2019). *High-dimensional statistics. A non-asymptotic viewpoint*, volume 48 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge.



# Sachverzeichnis

## A

a-posteriori-Risiko 14  
a-posteriori-Verteilung 14  
a-priori-Verteilung 13  
Ablehnbereich 21  
absolutstetig 138  
AIC 108  
Akaike-Informationskriterium 108  
Alternativhypothese 20  
asymptotisch normalverteilt 8

## B

balanciertes Design 87  
Bayes-Formel 139  
Bayes-Informationskriterium 116  
Bayes-optimal 14  
Bayes-Risiko 13, 15  
bedingte Dichte 139  
bedingte Wahrscheinlichkeit 132  
bedingter Erwartungswert 139  
Bestimmtheitsmaß 96  
Bias-Varianz-Zerlegung 7  
BIC 116  
Binomialtest  
  einseitiger 23, 36  
  zweiseitiger 25, 39  
BLUE 60  
Borel- $\sigma$ -Algebra 128  
Borel-Mengen 128

## C

charakteristische Funktion 136

## D

Design 49  
  zufällig 65  
Design-Kovarianzmatrix 65  
Designmatrix 53  
Dichte 129, 138  
dominiert 10  
Dummy-Variablen 53

## E

Effektdarstellung 92  
empirische Risiko 67  
empirischen Design-Kovarianzmatrix 68  
empirischer Risiko-Minimierer 67  
empirischer Risikominimierer  
  penalisierter 118  
empirisches Skalarprodukt 54, 123  
endlich 128  
erwartungstreu 6  
Erwartungswert 135

## F

F-Test 76, 82  
Faktor 87  
Fehler 1. und 2. Art 21  
Fehler- oder Störgrößen 53  
Fehlervektor 53  
Fisher-Test 76  
Fundamentalungleichung 118

## G

Gütefunktion 22  
Gauß-Test  
  einseitiger 33, 37  
  zweiseitiger 26, 39

**H**

Hypothesen 20

**I**

i.i.d. 134

**K**

Konfidenzintervall 40  
 Konfidenzmenge 40  
 konjugierte Verteilungsklassen 18  
 konsistent 8  
 Kontrastmatrix 80  
 Korrelation 136  
 Korrespondenzsatz 43  
 Kovarianz 136  
 Kovarianzmatrix 137  
 kritischer Bereich 21  
 kritischer Wert 25  
 Kullback-Leibler-Diskrepanz 106  
 Kullback-Leibler-Divergenz 103

**L**

Likelihood-Funktion 10  
 Likelihood-Quotiententest 38  
 lineares Modell 53  
   gewöhnliches 54

**M**

Maß 128  
   Dirac- 129  
   induziert 131  
   Lebesgue- 129  
   Zähl- 129  
 Maßraum 128  
 Mallows'  $C_p$ -Kriterium 111  
 mathematische Stichprobe 4  
 messbar 130  
 messbarer Raum 127  
 minimax 13  
 MLE 11  
 Modellmissspezifikation 102  
   Maximum-Likelihood-Schätzer unter 105  
 Momente 136  
 MSE 7

**N**

Neyman-Pearson-Lemma 35  
 Neyman-Pearson-Test 36

Nullhypothese 20

**O**

Orakelungleichung  
   Modellwahl 118  
 orthogonales Design 54  
 overfitting 101

**P**

p-Wert 30  
 Parametervektor 53  
 penalisierte Modellwahl 118

**Q**

Quantil-Quantil-Plot 27  
 Quantile 132  
 Quantilfunktion 132

**R**

Radon-Nikodym-Dichte 138  
 Regression 53  
   einfache 50  
   multiple 52  
   polynomiale 54  
 Regressionsgerade 50  
 Residuen 51, 60, 81  
 Responsevektor 53  
 Ridge-Regression 63  
 Risiko 7

**S**

$\sigma$ -Algebra 127  
 $\sigma$ -endlich 128  
 Satz  
   Gauß-Markov 58  
   Radon-Nikodym 138  
   von Bayes 133  
 Schätzer 6  
   Bayes- 14, 17  
   bester linearer erwartungstreuer 60  
   Kleinste-Quadrate- 55  
   Maximum-Likelihood- 11  
   Momenten- 9  
   Shrinkage- 63  
 signifikant 22  
 Signifikanzniveau 22  
 Standardabweichung 136  
 Statistik 4  
 statistischer Test  
   nichtrandomisierter 21  
   randomisierter 21  
 statistisches Experiment 3

- statistisches Modell 3
- Stichprobenmittel 6
- Stichprobenraum 3
- stochastisch beschränkt 141
- stochastische Ordnung 141
- Student-t-Test 78
- subgaußsch 68
- T**
- t-Test 78
  - Zweistichproben- 92
- Test zum Niveau  $\alpha$  22
- Testproblem 20
  - einseitiges 23
  - lineares 80
  - zweiseitig 23
- Teststatistik 24
- Totalvariationsabstand 122
- U**
- UMP-Tests 37
- Unabhängigkeit
  - von Ereignissen 133
  - von Zufallsvariablen 133
- underfitting 101
- Ungleichung
  - Cauchy-Schwarz 137
  - Hölder 137
  - Jensen 137
  - Markov 137
  - Tschebyscheff 137
- unkorreliert 136
- unverzerrt 6
- unverzerrte Risikoschätzung 110
- V**
- Varianz 136
- Varianzanalyse
  - einfaktorielle 87, 91
  - zweifaktorielle 93
- verallgemeinerte Inverse 131
- Verlust 7
  - 0-1- 7
  - absoluter 7
  - quadratischer 7
- Verteilung
  - $\chi^2$ - 145
  - Bernoulli- 142
  - Beta- 146
  - Binomial- 142
  - einer Zufallsvariable 131
  - F- 74
  - Gamma- 146
  - Gleich- 142, 144
  - hypergeometrische 143
  - inverse Gamma- 61
  - Multinomial- 142
  - Normal- bzw. Gauß- 144
  - normal-inverse Gamma- 62
  - Poisson- 144
  - t- 74
- Verteilungsfunktion 131
  - empirische 27
- Vorhersagefehler 66
- Vorhersageverlust 66
- W**
- Wahrscheinlichkeitsdichte 129
- Wahrscheinlichkeitsmaß 128
- Wahrscheinlichkeitsraum 128
- Wahrscheinlichkeitsverteilung
  - diskrete 129
  - stetige 129
- Wald-Intervall 42
- Wilson-Intervall 42, 48
- Z**
- Zähldichte 129
- zufälligem Design 65
- Zufallsvariablen 130