

# Mathematische Statistik

VON PETER PFAFFELHUBER

Version: 10. Februar 2016

**WARNUNG:** Dieses Skript enthält vermutlich noch viele Fehler. Es wurde teilweise in Eile geschrieben. Für alle Fehler bin ich selbst verantwortlich. Ich hoffe, in Zukunft eine Version mit weniger Fehlern bereit stellen zu können.

## Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>3</b>
1.1	Wiederholungen aus der Wahrscheinlichkeitstheorie . . . . .	3
1.2	Beispiele . . . . .	4
<b>2</b>	<b>Grundlegende Konzepte</b>	<b>6</b>
2.1	Suffizienz . . . . .	7
2.2	Vollständigkeit und Verteilungsfreiheit . . . . .	11
2.3	Exponentialfamilien . . . . .	14
2.4	Bayes'sche Modelle . . . . .	19
<b>3</b>	<b>Entscheidungstheorie</b>	<b>21</b>
3.1	Einführung . . . . .	21
3.2	Rolle suffizienter Statistiken . . . . .	26
3.3	Zulässige, Bayes, Minimax-Entscheidungsfunktionen . . . . .	28
<b>4</b>	<b>Testtheorie</b>	<b>36</b>
4.1	Bayes-Tests . . . . .	36
4.2	Likelihood-Quotienten-Tests . . . . .	37
4.3	Beste Tests . . . . .	42
<b>5</b>	<b>Schätztheorie</b>	<b>45</b>
5.1	Grundlagen . . . . .	46
5.2	UMVUE-Schätzer . . . . .	48
5.3	Information und die Cramér-Rao-Schranke . . . . .	50
5.4	Asymptotik von Maximum-Likelihood-Schätzern . . . . .	54

# 1 Einleitung

Die Mathematische Statistik ist ein eher theoretisches Teilgebiet der Stochastik. Im Gegensatz zur angewandten Statistik, die sich der Methodenentwicklung für Datenanalyse verschrieben hat, geht es in der theoretischen Statistik darum, Eigenschaften solcher Methoden festzustellen, etwa Optimalitätskriterien für Schätzer und Tests. In diesem Kurzschrift soll überblicksartig ein kurzer Einblick in dieses Gebiet gegeben werden. Komplettiert wird es (neben den Übungen) durch die Vorstellung statistischer Methoden, die an anderer Stelle im Rahmen dieser Vorlesung zusammengefasst vorgestellt werden.

## 1.1 Wiederholungen aus der Wahrscheinlichkeitstheorie

Wir setzen Kenntnisse aus den Vorlesungen *Stochastik I*, *Stochastik II* und *Wahrscheinlichkeitstheorie* voraus. Zur Sicherheit jedoch wiederholen wir einige Begriffe, die im Folgenden unerlässlich sein werden. Alle Räume in diesem Skript (z.B.  $E, E', E'', \mathbb{R}, \dots$ ) seien vollständige und separable metrische Räume, wenn nicht anders angegeben. Im Folgenden sei  $(E, \mathcal{A} = \mathcal{B}(E), \mathbb{P})$  ein Wahrscheinlichkeitsraum.

**Bemerkung 1.1 (Maß mit Dichte).** Das Bildmaß einer Zufallsvariable  $X$ , bezeichnet mit  $X_*\mathbb{P}$ , ist gegeben als  $X_*\mathbb{P}(A) = \mathbb{P}(X \in A)$ ,  $A \in \mathcal{B}(\mathbb{R})$ . Es hat Dichte  $p$  bezüglich  $\lambda^n$  (dem  $n$ -dimensionalen Lebesgue-Maß), falls für alle  $A \in \mathcal{B}(\mathbb{R})$

$$\mathbb{P}(X \in A) = \int 1_A(x)p(x)\lambda^n(dx).$$

In diesem Fall gilt dann für  $f : E \rightarrow \mathbb{R}$

$$\mathbb{E}[f(X)] = \int f(x)p(x)\lambda^n(dx),$$

falls eine der beiden Seiten existiert.

**Bemerkung 1.2 (Unabhängigkeit, Messbarkeit).** Zufallsvariablen  $X, Y$  sind (bezüglich  $\mathbb{P}$ ) unabhängig, wenn  $\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \cdot \mathbb{P}(Y \in B)$  für alle  $A, B \in \mathcal{B}(E)$ . Wir schreiben dann auch  $X \perp_{\mathbb{P}} Y$ .

Seien  $X, T$  Zufallsvariablen. Wir erinnern daran, dass  $\sigma(T) \subseteq \mathcal{A}$  die von  $T$  erzeugte  $\sigma$ -Algebra ist. Nach einem Satz aus der Wahrscheinlichkeitstheorie ist genau dann  $X$  messbar bezüglich  $\sigma(T)$  oder  $T$ -messbar, falls es eine Funktion  $g$  gibt mit  $X = g(T)$ .

**Bemerkung 1.3 (Bedingte Erwartung, bedingte Verteilung).** Sei  $(E, \mathcal{A}, \mathbb{P})$  ein Wahrscheinlichkeitsraum,  $X$  eine integrierbare Zufallsvariable und  $\mathcal{G} \subseteq \mathcal{F}$  eine (Teil-) $\sigma$ -Algebra. Die bedingte Erwartung von  $X$  gegeben  $\mathcal{G}$  (bezeichnet mit  $\mathbb{E}[X|\mathcal{G}]$ ) ist die einzige  $\mathcal{G}$ -messbare Zufallsvariable, für die

$$\mathbb{E}[X, G] = \mathbb{E}[\mathbb{E}[X|\mathcal{G}], G]$$

für alle  $G \in \mathcal{G}$  gilt.

Wir erinnern an zwei Eigenschaften der bedingten Erwartung: Für eine weitere  $\sigma$ -Algebra  $\mathcal{H} \subseteq \mathcal{G}$  gilt

$$\mathbb{E}[X|\mathcal{H}] = \mathbb{E}[\mathbb{E}[X|\mathcal{G}]|\mathcal{H}].$$

Für eine weitere reellwertige Zufallsvariable  $Y$  gilt

$$\mathbb{E}[X\mathbb{E}[Y|\mathcal{G}]] = \mathbb{E}[\mathbb{E}[X|\mathcal{G}]Y] = \mathbb{E}[\mathbb{E}[X|\mathcal{G}]\mathbb{E}[Y|\mathcal{G}]], \quad (1.1)$$

falls alle Erwartungen existieren.

Wir definieren weiterhin  $\mathbb{P}(A|\mathcal{H}) := \mathbb{E}[1_A|\mathcal{H}]$  für  $A \in \mathcal{F}$ . Angemerkt sei, dass es sich bei  $A \mapsto \mathbb{P}(A|\mathcal{H})$  zunächst nicht notwendigerweise um ein Wahrscheinlichkeitsmaß handelt, da bedingte Erwartungen zunächst nur  $\mathbb{P}$ -fast sicher definiert sind, es also immer Ausnahmeverteilungen geben kann. Ist jedoch  $E$  ein vollständiger und separabler metrischer Raum (was wir hier annehmen werden), so existiert (nach einem Satz aus der Wahrscheinlichkeitstheorie) die *reguläre Version der bedingten Verteilung*, d.h. ein stochastischer Kern  $\kappa$  von  $E$  nach  $E$ , so dass für  $\mathbb{P}$ -f.a.  $\omega \in E$

$$\kappa(\omega, B) = \mathbb{P}(X \in B|\mathcal{G})(\omega).$$

**Bemerkung 1.4 (Varianzzerlegung).** Neben der bedingten Erwartung kann man auch die bedingte Varianz

$$\mathbb{V}[X|\mathcal{G}] := \mathbb{E}[(X - \mathbb{E}[X|\mathcal{G}])^2|\mathcal{G}] = \mathbb{E}[X^2|\mathcal{G}] - \mathbb{E}[X|\mathcal{G}]^2$$

definieren. Es gilt dann die Varianzzerlegung

$$\mathbb{V}[X] = \mathbb{E}[\mathbb{V}[X|\mathcal{G}]] + \mathbb{V}[\mathbb{E}[X|\mathcal{G}]]. \quad (1.2)$$

**Bemerkung 1.5 (Notation).** Für  $x, y \in \mathbb{R}^n$  ist  $x^\top y$  das euklidische Skalarprodukt. Allgemeiner bezeichnen wir für  $A \in \mathbb{R}^{m \times n}$  und  $x \in \mathbb{R}^n$  die Matrixmultiplikation mit  $Ax$  und mit  $A^\top$  die Transponierte von  $A$ .

Wir bezeichnen das  $n$ -dimensionale Lebesgue-Maß (auf  $\mathbb{R}^n$ ) mit  $\lambda^n$ , und zur Vereinfachung der Notation  $n$ -dimensionale Zählmaß ebenfalls mit  $\lambda^n$  (also ist in diesem Fall  $\lambda^n = \sum_{x \in \mathbb{Z}^n} \delta_x$ ). Das  $n$ -dimensionale Produktmaß eines Maßes  $\mu$  bezeichnen wir im Allgemeinen mit  $\mu^n$ . Etwa ist für die Standardnormalverteilung  $\mathcal{N}(0, 1)$  die Verteilung einer unabhängigen Stichprobe gerade  $\mathcal{N}(0, 1)^n$ .

Für die vollständigen, separablen metrische Räume  $(D, r_D), (E, r_E), \dots$  seien  $\mathcal{B}(D), \mathcal{B}(E), \dots$  die Borel'schen  $\sigma$ -Algebren.

Hat  $\mu$  die Dichte  $f$  bezüglich  $\nu$ , so schreiben wir  $\mu = f \cdot \nu$ . Das Dirac-Maß auf  $x \in E$  bezeichnen wir mit  $\delta_x$ . Das Bildmaß von  $X$  unter  $\mu$  bezeichnen wir mit  $X_*\mu$ .

## 1.2 Beispiele

Als theoretische Wissenschaft mit Anwendungsbezügen lebt die Statistik von guten Beispielen, anhand denen man die zu entwickelnde Theorie ausprobieren kann. Auch wenn im Verlauf des Skriptes noch weitere Beispiele auftreten werden, sammeln wir hier drei besonders wichtige, die wir zunächst ohne großen Formalismus vorstellen.

**Beispiel 1.6 (Beispiel Bern).** Ein Bernoulli-Experiment besteht aus einer (endlich oder unendlich oft) unabhängig wiederholten Durchführung eines Zufallsexperiments, in dem jede Durchführung entweder einen *Erfolg* oder einen *Misserfolg* liefert. Es wird beschrieben durch einen Zufallsvektor  $X = (X_1, X_2, \dots)$  und eine Wahrscheinlichkeitsverteilung  $\mathbb{P}_\theta$ , so dass für  $x_i \in \{0, 1\}, i = 1, 2, \dots$

$$\mathbb{P}_\theta((X_1, \dots, X_n) = (x_1, \dots, x_n)) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} \quad (\text{Bern } 0)$$

für ein  $\theta \in [0, 1]$  gilt. Hierbei ist  $\theta$  die Wahrscheinlichkeit eines Erfolges (in jeder der Durchführungen). Bei einer gegebenen Folge von Erfolgen und Misserfolgen  $x_1, \dots, x_n$  wären naheliegende Fragen etwa:

- Wie groß ist  $\theta$ ?
- Ist  $\theta = \frac{1}{2}$ ?

**Beispiel 1.7 (Beispiel Norm).** Der zentrale Grenzwertsatz betont die Bedeutung der Normalverteilung. Für statistische Fragestellungen bedeutet dies, dass man – zumindest approximativ – oft eine Stichprobe von Daten  $X_1, \dots, X_n$  erhebt, die unabhängig und normalverteilt sind. Da man typischerweise die Stichprobe aus der gleichen Grundgesamtheit zieht, sollten hierbei die Erwartungswerte und Varianzen gleich sein. Hier ist also  $X = (X_1, \dots, X_n)$  und  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  hat für  $\theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+$  die Dichte

$$p_{(\mu, \sigma^2)}(x) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{\sigma^2}\right) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}\right) \quad (\text{Norm } 0)$$

bezüglich des Lebesgue-Maßes in  $\mathbb{R}^n$ . Bei einer Stichprobe  $X_1, \dots, X_n$  könnte man etwa fragen:

- Wie groß ist  $\mu$ , wie groß ist  $\sigma^2$ ?
- Ist  $\mu = \mu_0$  für einen vorgegebenen Wert  $\mu_0$ , wenn man weiß, wie groß  $\sigma^2$  ist?
- Ist  $\mu = \mu_0$  für einen vorgegebenen Wert  $\mu_0$ , wenn man nicht weiß, wie groß  $\sigma^2$  ist?

**Beispiel 1.8 (Beispiel Unif).** Etwas pathologischer klingt folgendes Beispiel: Wir nehmen an, dass Daten  $X = (X_1, \dots, X_n)$  unabhängig uniform auf  $[0, \theta]$  verteilt sind (für ein  $\theta \in \mathbb{R}_+$ ). Das heißt, dass  $\mathbb{P}_\theta$  so ist, dass  $X_*\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  mit

$$p_\theta(x) = \prod_{i=1}^n \frac{1}{\theta} \mathbb{1}_{0 \leq x_i \leq \theta} = \frac{1}{\theta^n} \mathbb{1}_{\max_{i=1, \dots, n} x_i \leq \theta}$$

für  $x \in \mathbb{R}_+^n$ . Hier könnte man etwa fragen:

- Wie groß ist  $\theta$ ?

Für alle Beispiele kann man also sagen, dass *Daten*  $X$  erhoben wurden, deren Verteilung von einem Parameter  $\theta$  abhängen. Dies führt zur ersten Definition.

**Definition 1.9 (Statistisches Modell).** Sei  $(\Omega, \mathcal{A})$  ein Messraum.

1. Ein statistisches Modell (auf  $E$ ) ist ein Paar  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$ , wobei  $X : \Omega \rightarrow E$  messbar und  $\{\mathbb{P}_\theta : \theta \in \mathcal{P}\}$  eine Familie von Wahrscheinlichkeitsmaßen auf  $\mathcal{A}$  ist. Hierbei heißt  $\mathcal{P}$  auch der Parameterraum.
2. Das statistische Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  heißt identifizierbar, wenn  $\theta = \theta'$  genau dann, wenn  $X_*\mathbb{P}_\theta = X_*\mathbb{P}_{\theta'}$  gilt. (Im weiteren Verlauf werden wir immer annehmen dass statistische Modelle diese Eigenschaft haben.)

3. Weiter heißt das statistische Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  regulär, falls  $E = \mathbb{R}^n$  für ein  $n$  und für alle  $\theta \in \mathcal{P}$

$$X_*\mathbb{P}_\theta = p_\theta \cdot \lambda^n$$

für eine geeignete Dichtefunktion  $p_\theta$  bzw. eine geeignete Zähldichte  $p_\theta$  gilt. (Wir erinnern daran, dass wir  $\lambda^n$  sowohl für das Lebesgue-Maß in  $\mathbb{R}^n$  oder das Zählmaß auf  $\mathbb{Z}^n$  verwenden.) Im ersten Fall sprechen wir von einem regulären stetigen Modell, im zweiten Fall von einem regulären, diskreten Modell.

**Bemerkung 1.10 (Parametrische und nicht-parametrische statistische Modelle).**

Ist  $\mathcal{P} \subseteq \mathbb{R}^k$  für ein  $k$ , so spricht man oft von parametrischen Modellen. Die Idee ist, dass  $\mathbb{P}_\theta$  von einem Vektor  $\theta$  von verschiedenen Parametern abhängt, etwa Mittelwert und Varianz einer Normalverteilung. In allen anderen Fällen spricht man von nicht-parametrischen Modellen. Im allgemeinsten Fall ist  $\mathcal{P}$  die Menge aller Wahrscheinlichkeitsmaße auf  $\mathcal{B}(\mathbb{R}^n)$ .

**Beispiel 1.11 (Beispiel Norm).** Im Fall von Beispiel 1.7 setzen wir  $\theta = (\mu, \sigma^2)$ , also

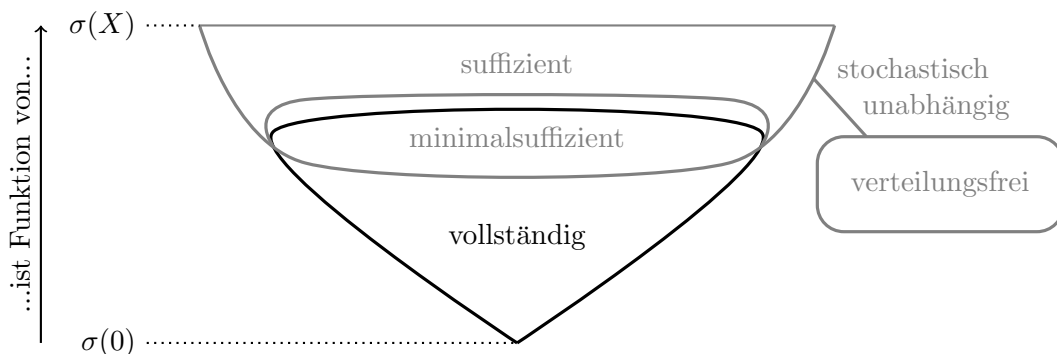
$$(X, \{\mathbb{P}_{\theta=(\mu, \sigma^2)} = \mathcal{N}(\mu, \sigma^2)^n : \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+\}). \quad (\text{Norm 1a})$$

Zu dieser Situation sagen wir auch, dass sowohl  $\mu$  als auch  $\sigma^2$  unbekannt sind. In einigen Situationen werden wir annehmen, dass wir etwa  $\sigma^2$  bereits kennen (aber  $\mu$  nicht). Dann setzen wir für dieses  $\sigma^2$  das statistische Modell

$$(X, \{\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2)^n : \theta \in \mathbb{R}\}). \quad (\text{Norm 1b})$$

## 2 Grundlegende Konzepte

In diesem und dem nächsten Abschnitt führen wir die Konzepte Suffizienz, Minimalsuffizienz, Vollständigkeit und Verteilungsfreiheit von Statistiken ein. (Eine Statistik ist dabei einfach eine  $\sigma(X)$ -messbare Zufallsvariable, d.h. für eine Abbildung  $t$  gilt  $T = t(X)$ .) Folgende Grafik illustriert kurz die Zusammenhänge.



Die wichtigsten Sätze des Abschnittes sind der Satz von Bahadur, Theorem 2.16, der besagt, dass suffiziente, vollständige Statistiken minimal-suffizient sind. Weiter besagt der Satz von Basu, Theorem 2.19, dass verteilungsfreie und suffiziente Statistiken unabhängig sind.

## 2.1 Suffizienz

Suffiziente Statistiken sind solche, die alle (zu statistischen Zwecken) nötigen Informationen über die Daten enthält.

**Definition 2.1 (Suffiziente Statistik).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell auf  $E$  und  $T = t(X)$  für  $t : E \rightarrow E'$  für einen vollständigen metrischen Raum  $(E', r')$  messbar. Dann heißt  $T$  suffizient, falls unter  $\mathbb{P}_\theta$  eine reguläre Version der bedingten Verteilung von  $X$  gegeben  $T$  existiert, die nicht von  $\theta$  abhängt. Insbesondere hängt also für  $A \in \mathcal{B}(E)$  die bedingte Erwartung  $\mathbb{P}_\theta(X \in A|T)$  nicht von  $\theta$  ab.

In jedem statistischen Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ist  $X$  suffizient. Weitere Beispiele für Suffizienz lassen sich mit Hilfe des Fisher-Neyman'schen Faktorisierungssatzes ausmachen, siehe Theorem 2.5.

**Bemerkung 2.2 (Umformulierung).** 1. Im Falle eines regulären, diskreten statistischen Modells bedeutet die bedingte Unabhängigkeit von  $X$  gegeben  $T = t(X)$  (unter  $\mathbb{P}_\theta$ ) gerade, dass (für  $t = t(x)$ )

$$\mathbb{P}_\theta(X = x|T = t) = \frac{\mathbb{P}_\theta(X = x)}{\mathbb{P}_\theta(t(X) = t)} = \frac{\mathbb{P}_\theta(X = x)}{\sum_{y:t(y)=t} \mathbb{P}_\theta(X = y)} = \frac{p_\theta(x)}{\sum_{y:t(y)=t} p_\theta(y)}$$

nicht von  $\theta$  abhängt. (Im Fall  $t \neq t(x)$  ist natürlich  $\mathbb{P}_\theta(X = x|t(X) = t) = 0$ .)

2. Die reguläre Version der bedingten Verteilung von  $X$  gegeben  $T$  ist (unter  $\mathbb{P}_\theta$ ) der  $\mathbb{P}_\theta$ -fast sicher eindeutige stochastische Kern  $\kappa_{X,T}$  (von  $E$  nach  $E'$ ), so dass

$$\kappa_{X,T,\theta}(\omega, A) = \mathbb{P}_\theta(X \in A|T)(\omega)$$

$\mathbb{P}_\theta$ -f.s. gilt; siehe auch Bemerkung 1.3. Diese existiert für vollständige und separable metrische Räume  $(E, r_E)$  nach einem Satz aus der Wahrscheinlichkeitstheorie. Diese reguläre Version der bedingten Verteilung ist dabei durch die Gleichung

$$\mathbb{E}_\theta[\mathbb{P}_\theta(X \in A|T), X \in B] = \mathbb{P}_\theta(X \in A \cap B)$$

für  $A \in \mathcal{B}(E)$  und  $B \in \sigma(T) = t^{-1}(\mathcal{B}(E'))$  eindeutig definiert. Genauer ist  $\mathbb{P}_\theta(X \in A|T)$  die  $\mathbb{P}_\theta$ -f.s. eindeutige,  $T$ -messbare Zufallsvariable, die diese Gleichung erfüllt.

**Beispiel 2.3 (Beispiel Bern).** Da das Beispiel 1.6 diskret ist, können wir hier nachrechnen, dass  $T := t(X) := \sum_{i=1}^n X_i$  suffizient ist. Beweis siehe Übung.

Zunächst erscheint es schwierig, suffiziente Statistiken auszumachen. Allerdings geben wir mit Theorem 2.5 eine einfache Charakterisierung von suffizienten Statistiken in Termen der Dichte von  $X$ . Erst einmal benötigen wir jedoch ein Lemma.

**Lemma 2.4 (Vorbereitung des Fisher-Neyman'schen Faktorisierungssatzes).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein reguläres, stetiges statistisches Modell mit  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$ , sowie  $t : E \rightarrow E'$  messbar und  $T = t(X)$ . Dann gilt:

1. Es gibt  $c_1, c_2, \dots \geq 0$  mit  $\sum_{i=1}^{\infty} c_i = 1$  und  $\theta_1, \theta_2, \dots \in \mathcal{P}$ , so dass  $\mathbb{P}_\theta \ll \nu^* := \sum_{i=1}^{\infty} c_i \mathbb{P}_{\theta_i}$  für alle  $\theta \in \mathcal{P}$ .

2. Ist  $T$  suffizient und  $A \in \mathcal{B}(E)$ , so gilt  $\nu^*(X \in A|T) \stackrel{\nu^*}{=} \mathbb{P}_\theta(X \in A|T)$  für alle  $\theta \in \mathcal{P}$ .

3. Gilt  $p_\theta(x) = h(x)g_\theta(t(x))$  für ( $\mathbb{P}_\theta$ -f.a.)  $x \in E, \theta \in \mathcal{P}$ , dann gibt es eine  $T$ -messbare Version von  $\frac{d\mathbb{P}_\theta}{d\nu^*}$ , nämlich

$$\frac{d\mathbb{P}_\theta}{d\nu^*} = \frac{g_\theta(t(x))}{\sum_{i=1}^{\infty} c_i g_{\theta_i}(t(x))}.$$

4. Falls  $\frac{d\mathbb{P}_\theta}{d\nu^*}$  für alle  $\theta \in \mathcal{P}$  nur von  $t(x)$  abhängt, so ist  $T = t(X)$  suffizient.

*Beweis.* 1. Sei  $P = p \cdot \lambda^n$  ein Wahrscheinlichkeitsmaß auf  $\mathbb{R}^n$  mit  $p > 0$  (etwa eine  $n$ -dimensionale, nicht-entartete Normalverteilung). Sei weiter

$$\mathcal{Q} = \left\{ \sum_{i=1}^{\infty} c_i \mathbb{P}_{\theta_i} \text{ für } \theta_1, \theta_2, \dots \in \mathcal{P} \text{ und } c_1, c_2, \dots \geq 0 \text{ mit } \sum_{i=1}^{\infty} c_i = 1 \right\}$$

die Menge der Konvexkombinationen aus  $\{\mathbb{P}_\theta : \theta \in \mathcal{P}\}$ . Wir setzen

$$\mathcal{C} := \left\{ C \in \mathcal{B}(E) : \exists Q \in \mathcal{Q} : Q(C) > 0 \text{ und } \frac{dQ}{dP}|_C > 0 \right\}.$$

Dann ist  $\mathcal{C}$  nicht leer und  $0 < \sup_{C \in \mathcal{C}} P(C) =: c \leq 1$ . Wähle nun  $C_1, C_2, \dots \in \mathcal{C}$  mit  $\sup_i P(C_i) = c$  und  $D := \bigcup_i C_i$  sowie  $Q_1, Q_2, \dots \in \mathcal{Q}$  so, dass  $Q_n(C_n) > 0$  und  $\frac{dQ_n}{dP}|_{C_n} > 0$ . Wir wählen

$$\nu^* = \sum_{i=1}^{\infty} 2^{-i} Q_i \in \mathcal{Q}.$$

Dann gilt sowohl  $\nu^*(D) > 0$  als auch  $d\nu^*/dP|_D = \sum_{i=1}^{\infty} 2^{-i} dQ_n/dP|_D > 0$  und damit  $D \in \mathcal{C}$ . Wir müssen zeigen, dass für  $\theta \in \mathcal{P}$  aus  $\nu^*(A) = 0$  folgt, dass  $\mathbb{P}_\theta(A) = 0$ . Wir setzen  $C := \{d\mathbb{P}_\theta/dP > 0\}$  und schreiben

$$\begin{aligned} \mathbb{P}_\theta(A) &= \mathbb{P}_\theta(A \cap D) + \mathbb{P}_\theta(A \cap D^c \cap C^c) + \mathbb{P}_\theta(A \cap D^c \cap C) \\ &\leq \int_{A \cap D} \frac{d\mathbb{P}_\theta}{dP} \cdot \left(\frac{d\nu^*}{dP}\right)^{-1} d\nu^* + \mathbb{P}_\theta(C^c) + \mathbb{P}_\theta(C \cap D^c). \end{aligned}$$

Die ersten beiden Summanden verschwinden. Angenommen,  $\mathbb{P}_\theta(C \cap D^c) > 0$ , dann wäre  $C \cup D \in \mathcal{C}$  und  $P(C \cup D) = P(D) + P(C \cap D^c) > P(D)$  im Widerspruch zur Maximalität von  $P(D)$ . Damit folgt  $\mathbb{P}_\theta(A) = 0$ .

2. Sei  $T$  suffizient. Für  $A \in \mathcal{B}(E)$  hängt  $\mathbb{P}_\theta(X \in A|T)$  ( $\mathbb{P}_\theta$ -f.s.) nicht von  $\theta$  ab, es gibt also eine Funktion  $\psi$  mit  $\mathbb{P}_\theta(X \in A|T) \stackrel{\mathbb{P}_\theta}{=} \psi(A, T)$  für alle  $\theta$ . Deshalb gilt auch

$$\nu^*(X \in A|T) = \sum_{i=1}^{\infty} c_i \mathbb{P}_{\theta_i}(X \in A|T) \stackrel{\nu^*}{=} \psi(A, T) \stackrel{\mathbb{P}_{\theta_i}}{=} \mathbb{P}_\theta(X \in A|T).$$

Die Behauptung folgt nun auch  $\nu^*$ -f.s., wenn man  $\psi$  auf  $\mathbb{P}_\theta$ -Nullmengen geeignet definiert.

3. Es gilt

$$\frac{d\nu^*}{d\lambda^n} = \sum_{i=1}^{\infty} c_i \frac{d\mathbb{P}_{\theta_i}}{d\lambda^n} = \sum_{i=1}^{\infty} c_i h(x) g_{\theta_i}(t(x))$$



und weiter, da  $\mathbb{P}_\theta \ll \nu^*$  für alle  $\theta \in \mathcal{P}$  (d.h. falls  $\nu^*(A) = 0$ , so ist zwar  $d\nu^*/d\lambda^n$  auf  $A$  nicht invertierbar, aber es gilt  $d\mathbb{P}_\theta/d\lambda^n = 0$ ) und

$$\frac{d\mathbb{P}_\theta}{d\nu^*} = \frac{d\mathbb{P}_\theta}{d\lambda^n} \left( \frac{d\nu^*}{d\lambda^n} \right)^{-1} = \frac{g_\theta(t(x))}{\sum_{i=1}^{\infty} c_i g_{\theta_i}(t(x))}.$$

Damit ist  $\frac{d\mathbb{P}_\theta}{d\nu^*}$  eine Funktion von  $t(x)$  und die Behauptung folgt.

4. Für  $A \in \mathcal{B}$ ,  $B \in \sigma(T)$  zeigen wir

$$\mathbb{E}_\theta[\mathbb{P}_\theta(X \in A|T), X \in B] = \mathbb{E}_\theta[\mathbb{P}_{\nu^*}(X \in A|T), X \in B].$$

Dann ist nämlich  $\mathbb{P}_\theta(X \in A|T)$  ( $\mathbb{P}_\theta$ -fs) unabhängig von  $\theta$  und  $T$  ist suffizient. Wir schreiben, da  $\frac{d\mathbb{P}_\theta}{d\nu^*}$  nach Voraussetzung messbar bezüglich  $\sigma(T)$  ist,

$$\begin{aligned} \mathbb{E}_\theta[\mathbb{P}_\theta(X \in A|T), X \in B] &= \mathbb{E}_\theta[1_{X \in A} 1_{X \in B}] = \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*}, X \in A \cap B \right] \\ &= \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*} \mathbb{P}_{\nu^*}(X \in A|T), X \in B \right] \\ &= \mathbb{E}_\theta[\mathbb{P}_{\nu^*}(X \in A|T), X \in B] \end{aligned}$$

und die Behauptung ist gezeigt.  $\square$

**Theorem 2.5 (Fisher-Neyman'scher Faktorisierungssatz).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell und  $T = t(X)$  mit  $t : \mathbb{R}^n \rightarrow \mathbb{R}^m$  messbar. Dann sind äquivalent:

1.  $T$  ist suffizient,
2. Es gibt  $g_\theta : \mathbb{R}^m \rightarrow \mathbb{R}$  und  $h : \mathbb{R}^n \rightarrow \mathbb{R}$ , so dass

$$p_\theta(x) = g_\theta(t(x))h(x).$$

**Bemerkung 2.6 (Diskreter Fall).** Im diskreten Fall ist der Beweis einfach. '2.  $\Rightarrow$  1.': Nach Bemerkung 2.2 gilt  $\mathbb{P}_\theta(X = x|t(X) = t) = 0$  für  $t \neq t(x)$ , was unabhängig von  $\theta$  ist. Für  $t = t(x)$  hingegen ist unter 2.

$$\mathbb{P}_\theta(X = x|t(X) = t) = \frac{g_\theta(t(x))h(x)}{\sum_{y:t(y)=t} g_\theta(t(y))h(y)} = \frac{g_\theta(t(x))h(x)}{g_\theta(t(x)) \sum_{y:t(y)=t} h(y)} = \frac{h(x)}{\sum_{y:t(y)=t} h(y)},$$

was ebenfalls unabhängig von  $\theta$  ist. Für '1.  $\Rightarrow$  2.' setzen wir

$$g_\theta(t) := \mathbb{P}_\theta(t(X) = t), \quad h(x) = \mathbb{P}_\theta(X = x|t(X) = t(x)).$$

Dann ist  $h(x)$  nach Voraussetzung unabhängig von  $\theta$  und es gilt

$$p_\theta(x) = \mathbb{P}_\theta(X = x) = \mathbb{P}_\theta(X = x, t(X) = t(x)) = h(x)g_\theta(t(x))$$

und die Behauptung ist gezeigt.

*Beweis im stetigen Fall.* '1.  $\Rightarrow$  2.': Sei  $\nu^*$  wie in Lemma 2.4.1. Nach Lemma 2.4.2 gilt  $\mathbb{P}_\theta(X \in A|T) = \nu^*(X \in A|T)$  ( $\nu^*$ -f.s.) für alle  $\theta \in \mathcal{P}$  und für  $A \in \mathcal{B}(E)$ . Wir schreiben nun mittels (1.1)

$$\begin{aligned} \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*}, X \in A \right] &= \mathbb{P}_\theta(X \in A) = \mathbb{E}_\theta[\mathbb{P}_\theta(X \in A|T)] \\ &= \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*} \nu^*(X \in A|T) \right] \\ &= \mathbb{E}_{\nu^*} \left[ \nu^*(X \in A|T) \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*} \middle| T \right] \right] \\ &= \mathbb{E}_{\nu^*} \left[ \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*} \middle| T \right], X \in A \right] \end{aligned}$$

wegen (1.1). Da  $A$  beliebig war und  $T = t(X)$ , gilt

$$\frac{d\mathbb{P}_\theta}{d\nu^*} = \mathbb{E}_{\nu^*} \left[ \frac{d\mathbb{P}_\theta}{d\nu^*} \middle| t(X) \right] =: g_\theta(t(X))$$

und mit  $h := \frac{d\nu^*}{d\lambda^n}$  gilt

$$p_\theta(x) = \frac{d\mathbb{P}_\theta}{d\nu^*} \frac{d\nu^*}{d\lambda^n} = g_\theta(t(x))h(x).$$

'2.  $\Rightarrow$  1.' Dies ist eine Folgerung aus Lemma 2.4.3 und Lemma 2.4.4.  $\square$

**Beispiel 2.7 (Beispiel Norm).** Bei normalverteilten Daten wie in Beispiel 1.7 und 1.11 schreiben wir für die Dichte

$$\begin{aligned} p_\theta(x) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{n\mu^2}{2\sigma^2}\right) \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2 + \frac{\mu}{\sigma^2} \sum_{i=1}^n x_i\right). \end{aligned}$$

Im Fall von unbekanntem  $\mu$  und  $\sigma^2$ , d.h. im statistischen Modell (Norm 1a) folgt mit dem Fisher-Neyman Kriterium, dass  $T = \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2\right)$  suffizient ist. Weiter können wir hier ablesen, dass im statistischen Modell (Norm 1b) (bei bekanntem  $\sigma^2$ ) bereits  $\sum_{i=1}^n X_i$  suffizient ist.

**Beispiel 2.8 (Beispiel Unif).** Für die Situation aus Beispiel 1.8 sehen wir mit dem Fisher-Neyman'schen Kriterium, dass  $T := t(X) := \max_{i=1, \dots, n} X_i$  suffizient ist. Beweis siehe Übung.

**Definition 2.9 (Minimalsuffizienz).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell und  $T = t(X)$  für  $t : E \rightarrow E'$  eine suffiziente Statistik. Dann heißt  $T$  minimal suffizient, falls es für jede weitere suffiziente Statistik  $U = u(X)$  für  $u : E \rightarrow E''$  eine messbare Abbildung  $g : E'' \rightarrow E'$  gibt mit  $T \stackrel{\mathbb{P}_\theta\text{-fs}}{=} g(U)$  für alle  $\theta \in \mathcal{P}$ .

Für obige Beispiele ist es nicht einfach nachzuprüfen, ob die angegebenen suffizienten Statistiken auch minimal suffizient sind. Folgender Satz erleichtert aber das Auffinden minimal suffizienter Statistiken.

**Theorem 2.10 (Kriterium für Minimalsuffizienz).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell mit  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  und  $t : E \rightarrow E'$  eine messbare Abbildung, so dass

$$t(y) = t(x) \iff \text{Es gibt ein } \ell(x, y) > 0, \text{ so dass für alle } \theta \in \mathcal{P}: p_\theta(y) = p_\theta(x)\ell(x, y).$$

Dann ist  $T = t(X)$  minimalsuffizient.

*Beweis.* Zunächst ist klar, dass die für  $t(y) = t(x)$  angegebene Bedingung symmetrisch in  $x$  und  $y$ , also wohldefiniert ist. (Mit  $\ell(x, y) > 0$  ist nämlich auch  $\ell(y, x) := 1/\ell(x, y) > 0$ .)

Wir zeigen zunächst, dass  $T$  suffizient ist. Sei  $\nu^*$  wie in Lemma 2.4. Es gilt für  $t(x) = t(y)$

$$\begin{aligned} \frac{d\mathbb{P}_\theta}{d\nu^*}(x) &= \frac{d\mathbb{P}_\theta}{d\lambda^n}(x) \Big/ \frac{d\nu^*}{d\lambda^n}(x) = \frac{p_\theta(x)}{\sum_{i=1}^{\infty} c_i p_{\theta_i}(x)} \\ &= \frac{p_\theta(x)\ell(x, y)}{\sum_{i=1}^{\infty} c_i p_{\theta_i}(x)\ell(x, y)} = \frac{p_\theta(y)}{\sum_{i=1}^{\infty} c_i p_{\theta_i}(y)} = \frac{d\mathbb{P}_\theta}{d\nu^*}(y). \end{aligned}$$

Damit hängt  $\frac{d\mathbb{P}_\theta}{d\nu^*}$  nur von  $t(x)$  ab und die Behauptung folgt nach Lemma 2.4.4.

Es folgt nun der Beweis, dass  $T$  minimalsuffizient ist. Sei hierzu  $U = u(X)$  eine weitere suffiziente Statistik. Nach dem Fisher-Neyman'schen Faktorisierungssatz gibt es  $g_\theta$  und  $h$ , so dass  $p_\theta(x) = g_\theta(u(x))h(x)$ . (Man kann oBdA annehmen, dass  $h > 0$  ist.) Wir müssen zeigen, dass aus  $u(x) = u(y)$  folgt, dass  $t(x) = t(y)$ . Dann nämlich gibt es eine Funktion  $g$  mit  $t(x) = g(u(x))$ . Sei also  $u(x) = u(y)$  und damit

$$\frac{p_\theta(y)}{p_\theta(x)} = \frac{g_\theta(u(y))h(y)}{g_\theta(u(x))h(x)} = \frac{h(y)}{h(x)}.$$

Daraus folgt  $t(x) = t(y)$  mit der Wahl  $\ell(x, y) = h(y)/h(x)$ .  $\square$

**Beispiel 2.11 (Beispiel *Bern*).** Wir haben bereits gesehen, dass  $T = \sum_{i=1}^n X_i$  suffizient ist. Außerdem gilt  $p_\theta(x) = p_\theta(y)$  genau dann, wenn  $t(x) = t(y)$ . Wählt man also  $\ell(x, y) = 1$  in Theorem 2.10, so erhält man die Minimalsuffizienz von  $T$ .

**Beispiel 2.12 (Beispiel *Unif*).** In Beispiel *Unif* ist  $t(x) := \max_{i=1, \dots, n} x_i$  suffizient; siehe Beispiel 2.8. Für festes  $\theta$  ist

$$\frac{p_\theta(y)}{p_\theta(x)} = \frac{\theta^n \mathbf{1}(t(y) \leq \theta)}{\theta^n \mathbf{1}(t(x) \leq \theta)} = \frac{\mathbf{1}(t(y) \leq \theta)}{\mathbf{1}(t(x) \leq \theta)}.$$

Nun  $t(x) = t(y)$  genau dann, wenn  $\frac{p_\theta(y)}{p_\theta(x)} = 1$  für alle  $\theta$ . Damit ist  $T = t(X)$  minimalsuffizient.

## 2.2 Vollständigkeit und Verteilungsfreiheit

Manchmal benötigt man suffiziente Statistiken, die weitere Eigenschaften erfüllen. Eine solche Eigenschaft wird nun beschrieben. Sie hilft insbesondere, minimalsuffiziente Statistiken zu finden; siehe Theorem 2.16.

**Definition 2.13 (Vollständige Statistik).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell. Eine Statistik  $T = t(X)$  für ein  $t : E \rightarrow E'$  heißt (beschränkt) vollständig, falls für alle (beschränkten) messbaren Funktionen  $g$  gilt, dass

$$\mathbb{E}_\theta[g(T)] = 0 \text{ für alle } \theta \in \mathcal{P} \implies g(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0 \text{ für alle } \theta \in \mathcal{P}.$$

**Beispiel 2.14 ( $X$  nicht beschränkt vollständig).** Wir zeigen nun, dass im Allgemeinen für ein statistisches Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  die Zufallsvariable  $X$  nicht vollständig ist. Sei hierzu etwa für ein  $n \geq 2$  die Verteilung  $\mathbb{P}_\theta = \text{Poi}(\theta)^n$  für  $\theta \in \mathcal{P} := \mathbb{R}_+$  die  $n$ -dimensionale Poisson-Verteilung mit Parameter  $\theta$ , d.h. unter  $\mathbb{P}_\theta$  hat  $X = (X_1, \dots, X_n)$  Werte in  $\mathbb{R}^n$  und  $X_1, \dots, X_n$  sind unabhängig und identisch Poisson verteilt zum Parameter  $\theta \geq 0$ . Für ein beschränktes (aber auf  $\mathbb{Z}_+$  nicht konstantes)  $f : \mathbb{Z}_+ \rightarrow \mathbb{R}$  sei  $g : \mathbb{Z}_+^n \rightarrow \mathbb{R}$  gegeben durch  $g(x_1, \dots, x_n) := f(x_1) - f(x_2)$ . Dann ist offensichtlich (da  $X_1 \sim X_2$  für alle  $\theta \in \mathcal{P}$ )

$$\mathbb{E}_\theta[g(X)] = \mathbb{E}_\theta[f(X_1)] - \mathbb{E}_\theta[f(X_2)] = 0,$$

aber  $\mathbb{P}_\theta(f(X_1) \neq f(X_2)) > 0$ , d.h.  $g(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0$  gilt nicht. Deshalb ist  $X$  nicht beschränkt vollständig (und damit nicht vollständig).

**Beispiel 2.15 (Poisson-verteilte Statistik  $T$ ).** Für ein statistisches Modell  $\{\mathbb{P}_\theta : \theta \in \mathbb{R}_+\}$  sei  $T_*\mathbb{P}_\theta = \text{Poi}(\theta)$ . Dann ist  $T$  beschränkt vollständig.

Denn: Sei  $g$  messbar und beschränkt, so dass

$$\mathbb{E}_\theta[g(T)] = e^{-\theta} \sum_{t=0}^{\infty} g(t) \frac{\theta^t}{t!} = 0$$

für alle  $\theta \geq 0$ . Dies ist die Potenzreihenentwicklung einer Funktion  $\theta \mapsto 0$ . Da diese Funktion analytisch ist, ist die Darstellung eindeutig und damit gilt  $g(t) = 0$  für  $t = 0, 1, 2, \dots$ , d.h.  $g(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0$ .

**Theorem 2.16 (Satz von Bahadur).** Sei  $\{\mathbb{P}_\theta : \theta \in \mathbb{R}_+\}$  ein statistisches Modell und  $T = t(X)$  für  $t : E \rightarrow \mathbb{R}^k$  für ein  $k \geq 0$  beschränkt vollständig und suffizient. Dann ist  $T$  minimalsuffizient.

*Beweis.* Wir setzen  $t(x) = (t_1(x), \dots, t_k(x))$ . Sei  $U = u(X)$  für  $u : E \rightarrow E''$  eine weitere suffiziente Statistik. Wir müssen zeigen, dass  $T$  eine Funktion von  $U$  ist. Wir setzen  $S = s(T)$  mit  $s = (s_1, \dots, s_k)$  und  $s_i(t) = (1 + e^{t_i})^{-1}$ . Dann ist  $s : \mathbb{R}^k \rightarrow \mathbb{R}^k$  injektiv und  $S$  ist beschränkt. Wir setzen für  $i = 1, \dots, k$

$$\begin{aligned} H_i(U) &= \mathbb{E}_\theta[S_i(T)|U], \\ L_i(T) &= \mathbb{E}_\theta[H_i(U)|T]. \end{aligned}$$

Da  $T$  und  $U$  suffizient sind, hängen diese Funktionen nicht von  $\theta$  ab und mit  $S$  sind auch  $H_i$  und  $L_i$  beschränkt,  $i = 1, \dots, k$ . Weiter gilt für  $\theta \in \mathcal{P}$

$$\mathbb{E}_\theta[S_i(T) - L_i(T)] = \mathbb{E}_\theta[H_i(U) - \mathbb{E}_\theta[H_i(U)|T]] = 0$$

und da  $T$  beschränkt vollständig ist folgt  $S_i(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} L_i(T)$  für alle  $\theta \in \mathcal{P}$ . Wir zerlegen nun die Varianz von  $L_i(T)$  und von  $H_i(U)$  mittels

$$\begin{aligned} \mathbb{V}_\theta[L_i(T)] &= \mathbb{E}_\theta[\mathbb{V}_\theta[L_i(T)|U]] + \mathbb{V}_\theta[\mathbb{E}_\theta[L_i(T)|U]] = \mathbb{E}_\theta[\mathbb{V}_\theta[L_i(T)|U]] + \mathbb{V}_\theta[H_i(U)], \\ \mathbb{V}_\theta[H_i(U)] &= \mathbb{E}_\theta[\mathbb{V}_\theta[H_i(U)|T]] + \mathbb{V}_\theta[L_i(T)] \end{aligned}$$

und damit  $\mathbb{V}_\theta[L_i(T)|U] = \mathbb{V}_\theta[H_i(U)|T] \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0$ , also auch  $\mathbb{V}_\theta[S_i(T)|U] \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0$ . Demnach gilt

$$S_i(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} \mathbb{E}_\theta[S_i(T)|U] \stackrel{\mathbb{P}_\theta\text{-fs}}{=} H_i(U).$$

Also gilt  $T_i = s_i^{-1}(H_i(U))$  für  $i = 1, \dots, k$  und  $T$  ist eine Funktion von  $U$ . Damit folgt die Behauptung.  $\square$

In Theorem 2.32 werden wir sehen, dass die suffizienten Statistiken aus Beispiel *Bern* und *Norm* vollständig sind. Beispiel *Unif* behandeln wir zunächst getrennt.

**Beispiel 2.17 (Beispiel *Unif*).** Wir haben bereits (siehe Beispiele 1.8, 2.8 und 2.12) gesehen, dass  $T := t(X) := \max_{i=1, \dots, n} X_i$  minimalsuffizient ist. Wir zeigen nun, dass  $T$  auch vollständig ist. Also würde hier zusammen mit der Suffizienz aus Beispiel 2.8 und dem Satz von Bahadur nochmals folgen, dass  $T$  minimalsuffizient ist.

Unter  $\mathbb{P}_\theta$  gilt

$$\mathbb{P}_\theta(T \leq t) = 1_{t \leq \theta} \left( \frac{t}{\theta} \right)^n,$$

also ist  $t \mapsto 1_{t \leq \theta} n t^{n-1} / \theta^n$  die Dichte von  $T$ . Gilt nun  $\mathbb{E}_\theta[g(T)] = 0$  für alle  $\theta > 0$ , so folgt

$$\int_0^\theta t^{n-1} g(t) dt = 0.$$

Da dies für alle  $\theta \geq 0$  gilt, ist

$$\int_a^b t^{n-1} g(t) dt = 0.$$

Dies ist nur möglich, wenn  $g \stackrel{\lambda\text{-f}}{=} 0$ .

Wir kommen nun zum Begriff der verteilungsfreien Statistik. Eine solche beinhaltet keine (statistisch verwendbaren) Informationen über die Daten. Das Hauptresultat über verteilungsfreie Statistiken ist Theorem 2.19, der den Zusammenhang zu suffizienten Statistiken herstellt; siehe auch die Abbildung am Anfang des Kapitels.

**Definition 2.18 (Verteilungsfreie Statistik).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell. Eine Statistik  $T = t(X)$  für ein  $t : E \rightarrow E'$  heißt verteilungsfrei, falls  $T_* \mathbb{P}_\theta$  unabhängig von  $\theta$  ist. Eine verteilungsfreie Statistik  $T$  heißt maximal, wenn es für jede andere verteilungsfreie Statistik  $U$  eine Funktion  $g$  gibt mit  $U = g(T)$ .

**Theorem 2.19 (Satz von Basu).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell. Weiter sei  $T = t(X)$  für  $t : E \rightarrow E'$  und  $U = u(X)$  für  $u : E \rightarrow E''$ .<sup>1</sup>

1. Ist  $T$  beschränkt vollständig und suffizient sowie  $U$  verteilungsfrei, dann gilt  $T \perp_{\mathbb{P}_\theta} U$  für alle  $\theta \in \mathcal{P}$ .
2. Angenommen, für alle  $\theta, \theta' \in \mathcal{P}$  gibt es eine Menge  $A \in \mathcal{B}(E)$  mit  $\mathbb{P}_\theta(X \in A), \mathbb{P}_{\theta'}(X \in A) > 0$ . Ist  $T \perp_{\mathbb{P}_\theta} U$  für alle  $\theta \in \mathcal{P}$  und  $T$  suffizient, dann ist  $U$  verteilungsfrei.
3. Sei  $T \perp_{\mathbb{P}_\theta} U$  für alle  $\theta \in \mathcal{P}$  und  $U$  verteilungsfrei. Wenn  $\sigma(T, U) = \sigma(X)$ , dann ist  $T$  suffizient.

*Beweis.* 1. Sei  $A \in \mathcal{B}(E'')$ . Offenbar gilt

$$\mathbb{P}_\theta(U \in A) = \mathbb{E}_\theta[\mathbb{P}_\theta(U \in A|T)].$$

Weiter hängt (wegen der Verteilungsfreiheit von  $U$ ) weder  $\mathbb{P}_\theta(U \in A)$  noch (wegen der Suffizienz von  $T$ ) die Größe  $\mathbb{P}_\theta(U \in A|T)$  von  $\theta$  ab. Damit ist  $g : t \mapsto \mathbb{P}_\theta(U \in A) - \mathbb{P}_\theta(U \in A|T)$  eine beschränkte messbare Funktion (unabhängig von  $\theta$ ) mit  $\mathbb{E}_\theta[g(T)] = 0$  für alle  $\theta \in \mathcal{P}$ .

<sup>1</sup>Wir erinnern an die Schreibweise  $X \perp_{\mathbb{P}_\theta} Y$ , falls  $X$  und  $Y$  unter  $\mathbb{P}_\theta$  stochastisch unabhängig sind.

Wegen der beschränkten Vollständigkeit von  $T$  ist damit  $\mathbb{P}_\theta(U \in A) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} \mathbb{P}_\theta(U \in A|T)$ . Dies aber bedeutet die Unabhängigkeit von  $U$  und  $T$ .

2. Da  $T$  suffizient ist, hängt  $\mathbb{P}_\theta(U \in A|T)$  für alle  $A \in \mathcal{B}(E'')$  nicht von  $\theta$  ab. Da  $\mathbb{P}_\theta(U \in A|T) = \mathbb{P}_\theta(U \in A)$  wegen der Unabhängigkeit von  $U$  und  $T$ , hängt also  $\mathbb{P}_\theta(U \in A)$  nicht von  $\theta$  ab. Dies ist aber gerade die Verteilungsfreiheit von  $U$ , da  $A$  beliebig war.

3. Es ist zu zeigen, dass für alle  $A \in \mathcal{B}(E)$  gilt, dass  $\mathbb{P}_\theta(X^{-1}(A)|T)$  unabhängig von  $\theta$  ist. In der Tat genügt es mittels einer monotonen Klasse, für einen schnittstabilen Erzeuger  $\mathcal{E}$  von  $\sigma(X)$  zu zeigen, dass für alle  $E \in \mathcal{E}$   $\mathbb{P}_\theta(E|T)$  unabhängig von  $\theta$  ist.

Sei nun  $B \in \mathcal{B}(E')$  und  $C \in \mathcal{B}(E'')$ . Es gilt

$$\mathbb{P}_\theta(T \in B, U \in C|T) = 1_{T \in B} \mathbb{P}_\theta(U \in C)$$

und dies ist (wegen der Verteilungsfreiheit von  $U$ ) unabhängig von  $\theta$ . Da  $\{T^{-1}(B) \cap U^{-1}(C) : B \in \mathcal{B}(E'), C \in \mathcal{B}(E'')\}$  schnittstabil ist und nach Voraussetzung  $\sigma(X)$  erzeugt, folgt die Aussage.  $\square$

**Beispiel 2.20 (Beispiel Norm).** Für festes  $\sigma^2$  betrachten wir das statistische Modell aus (Norm 1b), also das Normalverteilungsmodell mit bekanntem  $\sigma^2$ . Weiter setzen wir

$$T := \bar{X} := \frac{1}{n} \sum_{i=1}^n X_i, \quad U = \frac{1}{n-1} \sum_{i=1}^n (X_i - T)^2.$$

Aus der Rechnung in Beispiel 2.7 sieht man zusammen mit dem Fisher-Neyman'schen Faktorisierungssatzes (Theorem 2.5), dass  $T$  suffizient ist. (Weiter werden wir in Theorem 2.32 sehen, dass  $T$  auch vollständig, also sogar minimalsuffizient ist.) Außerdem ist  $U$  verteilungsfrei, denn: Der Vektor  $(X_1 - T, \dots, X_n - T)$  ist unabhängig von  $\theta$  normalverteilt mit

$$\begin{aligned} \mathbb{E}_\theta[X_i - T] &= 0, \\ \text{COV}_\theta[X_i - T, X_j - T] &= \sigma^2(\delta_{ij} - \frac{2}{n} + \frac{1}{n}) = \sigma^2(\delta_{ij} - \frac{1}{n}). \end{aligned}$$

Nun ist  $U$  eine Funktion dieses Vektors, also verteilungsfrei.) Nach dem Theorem von Basu, Theorem 2.19.2, sind diese beiden Vektoren also unabhängig. (Dieses Ergebnis ist auch Teil des bekannten Satzes von Fisher.)

## 2.3 Exponentialfamilien

Viele Verteilungen, etwa die Normal-, Poisson-, Binomial- und Exponentialverteilung, haben eine gemeinsame Struktur, die oftmals direkte Rechnungen ermöglicht. Diese Struktur wird in der folgenden Definition formalisiert.

**Definition 2.21 (Exponentialfamilie).** Sei  $\mathcal{P} \subseteq \mathbb{R}^k$ . Eine Familie  $\{\mathbb{P}_\theta : \theta \in \mathcal{P}\}$  von Wahrscheinlichkeitsmaßen auf  $\mathbb{R}^n$  (auf  $\mathbb{Z}^n$ ) heißt  $k$ -parametrische Exponentialfamilie (mit  $c, t, d, h$ ) für

$$c_1, \dots, c_k, d : \mathcal{P} \rightarrow \mathbb{R}, \quad t_1, \dots, t_k, h : \mathbb{R}^n \rightarrow \mathbb{R},$$

falls  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  und

$$p_\theta(x) = h(x) \cdot \exp\left(\sum_{j=1}^k c_j(\theta) t_j(x) - d(\theta)\right) = h(x) \cdot \exp(c(\theta)^\top t(x) - d(\theta)), \quad x \in \mathbb{R}^n.$$

Gilt insbesondere  $c_j(\theta) = \theta_j$ , also

$$p_\theta(x) = h(x) \cdot \exp(\theta^\top t(x) - d(\theta)), \quad x \in \mathbb{R}^n,$$

so sagt man, die Exponentialfamilie sei in kanonischer Form und wir definieren den kanonischen Parameterraum

$$\Gamma := \left\{ \theta \in \mathbb{R}^k : \int h(x) e^{\theta^\top t(x)} d\lambda^n(x) < \infty \right\}.$$

**Bemerkung 2.22 (1-parametrische Exponentialfamilie).** Für den Spezialfall einer 1-parametrischen Exponentialfamilie gibt es Funktionen  $c, d, t, h$  mit

$$p_\theta(x) = h(x) \cdot \exp(c(\theta)t(x) - d(\theta)), \quad x \in \mathbb{R}^n.$$

**Beispiel 2.23 (Beispiel Norm).** Wir betrachten das statistische Modell (Norm 1a) für  $n = 1$ , also  $\theta := (\mu, \sigma^2)$  mit  $\mathcal{P} = \mathbb{R} \times \mathbb{R}_+$  und  $\mathbb{P}_\theta := \mathcal{N}(\mu, \sigma^2)$ , und damit

$$p_{(\mu, \sigma^2)}(x) = \exp\left(\frac{\mu}{\sigma^2}x - \frac{x^2}{2\sigma^2} - \frac{1}{2}\left(\frac{\mu}{\sigma^2} + \log(2\pi\sigma^2)\right)\right).$$

Also ist die Familie der (ein-dimensionalen) Normalverteilungen eine 2-parametrische Exponentialfamilie mit

$$\begin{aligned} c_1(\mu, \sigma^2) &= \frac{\mu}{\sigma^2}, & t_1(x) &= x, \\ c_2(\mu, \sigma^2) &= -\frac{1}{2\sigma^2}, & t_2(x) &= x^2, \\ h(x) &= 1, & d(\mu, \sigma^2) &= -\frac{1}{2}\left(\frac{\mu}{\sigma^2} + \log(2\pi\sigma^2)\right). \end{aligned}$$

Diese ist nun allerdings nicht in kanonischer Form.

**Beispiel 2.24 (Beispiel Bern).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in (0, 1)\})$  wie in Beispiel 1.6. Dann ist also mit (Bern 0)

$$\begin{aligned} p_\theta(x) &= \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} \\ &= \exp\left(\sum_{i=1}^n x_i \log \theta + \left(n - \sum_{i=1}^n x_i\right) \log(1 - \theta)\right) \\ &= \exp\left(\log \frac{\theta}{1 - \theta} \sum_{i=1}^n x_i + n \log(1 - \theta)\right) \end{aligned}$$

und damit ist  $\{\mathbb{P}_\theta : \theta \in (0, 1)\}$  eine 1-parametrische Exponentialfamilie.

**Beispiel 2.25 (Beispiel Unif).** Die Familie  $\{\mathbb{P}_\theta : \theta \in \mathbb{R}_+\}$  aus Beispiel 1.8 ist keine Exponentialfamilie.

**Beispiel 2.26 (Beispiel Exp).** Sei  $\mathcal{P} = \mathbb{R}_+$  und  $\mathbb{P}_\theta = p_\theta \cdot \lambda$  die Exponentialverteilung mit Parameter  $\theta$ . Dann ist

$$p_\theta(x) = 1_{x \geq 0} e^{-\theta x + \log \theta}$$

und damit ist die Familie der Exponentialverteilungen eine 1-parametrische Exponentialfamilie und obige Darstellung ist die kanonische Form.

**Beispiel 2.27 (Beispiel *Pois*).** Sei  $\Lambda = \mathbb{R}_+$ . Für die Poisson-Verteilung mit Parameter  $\lambda$  schreiben wir  $\mathbb{P}_\lambda = p_\lambda \cdot \mu$  mit

$$p_\lambda(x) = \frac{1_{x \geq 0}}{x!} \exp((\log \lambda)x - \lambda)$$

Diese ist damit eine 1-parametrische Exponentialfamilie. Um obige Darstellung in kanonische Form zu bringen, setzen wir  $\theta := \log \lambda$  und schreiben für  $\theta \in \mathbb{R}$

$$p_\theta(x) = \frac{1_{x \geq 0}}{x!} \exp(\theta x - e^\theta).$$

**Bemerkung 2.28 (Sample aus einer Exponentialfamilie).** Ist  $\{\mathbb{P}_\theta : \theta \in \mathcal{P}\}$  eine  $k$ -parametrische Exponentialfamilie auf  $\mathbb{R}^n$  mit Funktionen  $c_1, \dots, c_k, d, t_1, \dots, t_k, h$ , und sind  $X_1, \dots, X_N$  unabhängig und identisch nach  $\mathbb{P}_\theta$  verteilt. Dann ist die gemeinsame Verteilung von  $X_1, \dots, X_N$  eine Exponentialfamilie mit

$$c_1, \dots, c_k, Nd, \sum_{i=1}^N t_1 \circ \pi_i, \dots, \sum_{i=1}^N t_k \circ \pi_i, \prod_{i=1}^N h \circ \pi_i$$

mit der Projektion  $\pi_i(x) = x_i$ .

Denn: Als gemeinsame Dichte schreiben wir

$$p_\theta(x_1, \dots, x_N) = h(x_1) \cdots h(x_N) \cdot \exp\left(\sum_{j=1}^k c_j(\theta) \sum_{i=1}^N t_j(x_i) - Nd(\theta)\right).$$

**Proposition 2.29 (Suffiziente Statistik bei Exponentialfamilien).**

Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine Exponentialfamilie (mit  $c, t, d, h$ ). Dann ist die Statistik  $T := t(X) = (t_1(X), \dots, t_k(X))$  suffizient.

*Beweis.* Um die Suffizienz von  $T$  zu sehen, schreiben wir zunächst

$$p_\theta(x) = h(x)g_\theta(t(x))$$

für

$$g_\theta(t(x)) = \exp(c(\theta)^\top t(x) - d(\theta)).$$

Damit folgt die Aussage aus dem Fisher-Neyman'schen Faktorisierungssatz, Theorem 2.5.  $\square$

**Beispiel 2.30 (Lineare Regression).** Bei der linearen Regression beobachten wir  $(x_1, Y_1), \dots, (x_n, Y_n)$  (mit  $x_i = (x_{i0} = 1, x_{i1}, \dots, x_{im}) \in \mathbb{R}^{m+1}$ ) und gehen davon aus, dass

$$Y_i = x_i \beta + \varepsilon_i \quad (\text{also } Y = x^\top \beta + \varepsilon)$$

für  $\beta_0, \dots, \beta_m \in \mathbb{R}$  mit  $\varepsilon_1, \dots, \varepsilon_n$  unabhängig und identisch nach  $\mathcal{N}(0, \sigma^2)$  verteilt. Hierbei sind  $x_1, \dots, x_n$  bekannt (und fest, also keine Parameter des Modells) und  $Y_1, \dots, Y_n$  werden als



Zufallsvariablen betrachtet. Wir haben also ein statistisches Modell  $(Y, \{\mathbb{P}_\theta = p_\theta \cdot \lambda^n : \theta = (\beta, \sigma^2) \in \mathbb{R}^{m+2}\})$  und wir schreiben für die gemeinsame Dichte von  $Y = (Y_1, \dots, Y_n)$

$$\begin{aligned} p_\theta(y) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i\beta)^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 + \frac{1}{\sigma^2} \beta^\top \sum_{i=1}^n y_i x_i - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i\beta)^2 - \frac{n}{2} \log(2\pi\sigma^2)\right) \end{aligned}$$

Wählen wir

$$\begin{aligned} c_j(\beta) &= \frac{\beta_j}{\sigma^2}, & t_j(y) &= \sum_{i=1}^n y_i x_{ij}, & j &= 0, \dots, m, \\ c_{m+1}(\beta) &= -\frac{1}{2\sigma^2}, & t_{m+1}(y) &= \sum_{i=1}^n y_i^2, \end{aligned}$$

$$d(\beta) = \frac{1}{2\sigma^2} \sum_{i=1}^n (\beta^\top x_i)^2 + \frac{n}{2} \log(2\pi\sigma^2),$$

so sehen wir, dass  $\{\mathbb{P}_\theta : \theta \in \mathbb{R}^{m+2}\}$  eine  $m+2$ -parametrische Exponentialfamilie ist. Weiter ist

$$\left( \sum_{i=1}^n Y_i, \sum_{i=1}^n Y_i x_{i1}, \dots, \sum_{i=1}^n Y_i x_{im}, \sum_{i=1}^n Y_i^2 \right)$$

suffizient.

**Lemma 2.31 (Kanonischer Parameterraum konvex).** *Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine Exponentialfamilie in kanonischer Form (mit  $c, t, d, h$ ). Dann ist der kanonische Parameterraum konvex und  $\theta \mapsto e^{d(\theta)}$  ist konvex.*

*Beweis.* Sei  $\theta, \eta \in \Gamma$  und  $0 < \alpha < 1$ . Dann gilt, da die Exponentialfunktion konvex ist,

$$\begin{aligned} e^{d(\alpha\theta + (1-\alpha)\eta)} &:= \int h(x) \exp((\alpha\theta + (1-\alpha)\eta)^\top t(x)) d\lambda^n(x) \\ &\leq \int h(x) (\alpha e^{\theta^\top t(x)} + (1-\alpha) e^{\eta^\top t(x)}) d\lambda^n(x) \\ &= \alpha e^{d(\theta)} + (1-\alpha) e^{d(\eta)}. \end{aligned}$$

Dies zeigt alle Behauptungen. □

**Theorem 2.32 (Vollständigkeit der suffizienten Statistik).** *Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine  $k$ -parametrische Exponentialfamilie in kanonischer Form (mit  $(\theta, t, d, h)$ ). Ist  $\mathbb{R}^k \supseteq \mathcal{P}^\circ \neq \emptyset$ , so ist  $T = t(X)$  vollständig (und wegen der Suffizienz aus Proposition 2.29 mit Theorem 2.16 auch minimal-suffizient).*

<sup>2</sup>Mit  $A^\circ$  bezeichnen wir das Innere der Menge  $A$ .

*Beweis.* Wir beweisen die Behauptung nur im Fall  $k = 1$ , da die anderen Fälle mit Induktion folgen. Angenommen,  $T$  sei nicht vollständig. Dann gibt es  $g : \mathbb{R} \rightarrow \mathbb{R}$  messbar und  $\mathbb{E}_\theta[g(T)] = 0$  für alle  $\theta \in \mathcal{P}$ , aber  $\mathbb{P}_{\theta_0}(g(T) > 0) > 0$  für ein  $\theta_0 \in \mathcal{P}^\circ$ . (Wir können oBdA  $\theta_0$  im Inneren von  $\mathcal{P}$  wegen der Stetigkeit der Dichte von  $T$  und monotoner Konvergenz wählen.) Anders ausgedrückt heißt das, dass für alle  $\theta \in \mathcal{P}$

$$\mathbb{E}_\theta[g^+(T)] = \mathbb{E}_\theta[g^-(T)] \quad (*)$$

gilt, aber

$$\mathbb{E}_{\theta_0}[g^+(T)] = \mathbb{E}_{\theta_0}[g^-(T)] =: w \in (0, \infty)$$

für ein  $\theta_0 \in \mathcal{P}^\circ$ . Wir definieren die beiden Wahrscheinlichkeitsmaße

$$P := \frac{1}{w}g^+ \circ T \cdot \mathbb{P}_{\theta_0}, \quad Q := \frac{1}{w}g^- \circ T \cdot \mathbb{P}_{\theta_0}.$$

Nun schreiben wir (\*) um in

$$\begin{aligned} \mathbb{E}_P[e^{(\theta-\theta_0)T}] &= \frac{1}{w}\mathbb{E}_{\theta_0}[e^{(\theta-\theta_0)T}g^+(T)] \\ &= \frac{1}{w} \int e^{(\theta-\theta_0)t(x)}g^+(t(x))h(x)e^{\theta_0t(x)-d(\theta_0)}d\lambda(x) \\ &= \frac{e^{d(\theta)-d(\theta_0)}}{w} \int g^+(t(x))h(x)e^{\theta t(x)-d(\theta)}d\lambda(x) \\ &= \frac{e^{d(\theta)-d(\theta_0)}}{w}\mathbb{E}_\theta[g^+(T)] = \frac{e^{d(\theta)-d(\theta_0)}}{w}\mathbb{E}_\theta[g^-(T)] \\ &= \frac{1}{w}\mathbb{E}_{\theta_0}[e^{(\theta-\theta_0)T}g^-(T)] \\ &= \mathbb{E}_Q[e^{(\theta-\theta_0)T}]. \end{aligned}$$

Wegen der Eindeutigkeit der Laplace-Transformierten bedeutet dies  $P = Q$  und damit auch  $g^+ \stackrel{\mathbb{P}_{\theta_0}\text{-fs}}{=} g^-$ , also  $g \stackrel{\mathbb{P}_{\theta_0}\text{-fs}}{=} 0$  und damit (wegen der Form der Exponentialfamilie) auch  $g \stackrel{\mathbb{P}_\theta\text{-fs}}{=} 0$  für alle  $\theta \in \mathcal{P}$ .  $\square$

**Bemerkung 2.33 (Exponentialfamilien in nicht-kanonischer Form).** Das Ergebnis aus Theorem 2.32 lässt sich durch Umparametrisierung auf Exponentialfamilien übertragen, die nicht in kanonischer Form vorliegen. Man sieht etwa, dass die suffizienten Statistiken aus Beispiel 2.23, 2.24, 2.26 und 2.27 minimalsuffizient sind.

Exponentialfamilien erlauben oft explizite Berechnungen. Dies illustrieren wir an zwei Ergebnissen, von denen wir das erste ohne Beweis angeben.

**Lemma 2.34 (Differenzierbarkeit nach  $\theta$ ).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine Exponentialfamilie in kanonischer Form (mit  $\theta, t, d, h$ ) und der kanonische Parameterraum  $\Gamma$  habe nicht-leeres Inneres,  $\Gamma^\circ \neq \emptyset$ . Weiter sei  $\phi : \mathbb{R}^k \rightarrow \mathbb{R}$  messbar und  $\theta \in \Gamma^\circ$  so, dass  $\eta \mapsto \mathbb{E}_\eta[|\phi(t(X))|e^{\eta^\top t(X)}] < \infty$  für  $\eta$  in einer Umgebung von  $\theta$ . Dann ist die Abbildung

$$f : \eta \mapsto \mathbb{E}_\eta[\phi(t(X))e^{\eta^\top t(X)}]$$

in einer Umgebung von  $\theta$  analytisch mit

$$\frac{\partial f(\eta)}{\partial \eta_i} = \mathbb{E}_\eta[t_i(X)\phi(t(X))e^{\eta^\top t(X)}].$$

**Proposition 2.35 (Laplace-Transformierte von Exponentialfamilien).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine  $k$ -parametrische Exponentialfamilie in kanonischer Form. Ist  $\theta \in \mathcal{P}^\circ$ , dann existiert die Laplace-Transformierte von  $t(X)$  mit

$$\mathbb{E}_\theta[e^{\eta^\top t(X)}] = \exp(d(\eta + \theta) - d(\theta)),$$

falls  $|\eta|$  klein genug ist und für  $l_1, \dots, l_k \geq 0$  mit  $l_1 + \dots + l_k = l$

$$\mathbb{E}_\theta \left[ \prod_{i=1}^k t_i(X)^{l_i} \right] = e^{-d(\theta)} \frac{\partial^l}{\partial \eta_1^{l_1} \dots \partial \eta_k^{l_k}} e^{d(\eta)} \Big|_{\eta=\theta}.$$

Insbesondere gilt also

$$\begin{aligned} \mathbb{E}_\theta[t_i(X)] &= \frac{\partial d(\eta)}{\partial \eta_i} \Big|_{\eta=\theta}, \\ \text{COV}_\theta[t_i(X), t_j(X)] &= \frac{\partial^2 d(\eta)}{\partial \eta_i \partial \eta_j} \Big|_{\eta=\theta}. \end{aligned}$$

*Beweis.* Wir berechnen

$$\begin{aligned} \mathbb{E}[e^{\eta^\top t(X)}] &= \int \exp(\eta^\top t(x)) \cdot h(x) \cdot \exp(\theta^\top t(x) - d(\theta)) \lambda^n(dx) \\ &= e^{d(\eta+\theta) - d(\theta)} \int h(x) \cdot \exp((\theta + \eta)^\top t(x) - d(\theta + \eta)) \lambda^n(dx) \\ &= e^{d(\eta+\theta) - d(\theta)}, \end{aligned}$$

da das Integral über eine Dichte eins ist. Die zweite Behauptung folgt aus Lemma 2.34.  $\square$

## 2.4 Bayes'sche Modelle

Die Formel von Bayes ist wohlbekannt. Auf ihr beruht der große Zweig der Bayesianischen Statistik. Grundlegend ist hier, dass in einem statistischen Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein Vorwissen über die Möglichkeiten besteht, welcher Parameter  $\theta \in \mathcal{P}$  zutrifft. Dies wird in der a-priori-Verteilung zusammengefasst, einer Verteilung auf  $\mathcal{P}$ .

**Definition 2.36 (A-priori-Verteilung, a-posteriori-Verteilung).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell und  $\mathcal{P}$  ein vollständiger, separabler metrischer Raum. Eine a-priori-Verteilung ist eine Wahrscheinlichkeitsverteilung  $\pi$  auf  $\mathcal{P}$ . Ist für alle  $A \in \mathcal{B}(E)$  die Abbildung  $\theta \mapsto \mathbb{P}_\theta(X \in A)$  messbar, so wird durch  $(\theta, A) \mapsto \mathbb{P}_\theta(X \in A)$  ein Markov-Kern von  $\mathcal{P}$  nach  $\mathcal{B}(E)$  definiert und für  $\Theta \sim \pi$  wird durch

$$P(\Theta \in A, X \in B) := \int_A \pi(d\theta) \mathbb{P}_\theta(X \in B)$$

die gemeinsame Verteilung von  $\Theta$  und  $X$  auf  $\mathcal{B}(\mathcal{P}) \otimes \mathcal{B}(E)$  definiert. Die a-posteriori-Verteilung ist dann die reguläre Version der bedingten Verteilung  $P(\Theta \in \cdot | X)$ , also

$$\pi_x = P(\Theta \in \cdot | X = x).$$

**Bemerkung 2.37 (A-posteriori-Verteilung bei regulären Modellen).** Ist  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell und  $\mathcal{P} \subseteq \mathbb{R}^k$ , und ist  $\pi = g \cdot \lambda^k$ , so hat die gemeinsame Verteilung von  $\Theta$  und  $X$  die Dichte  $g(d\theta) \cdot p_\theta(dx)$ . Die a-posteriori Verteilung  $\pi_x$  hat dann die Dichte

$$p_x(\theta) = \frac{p_\theta(x)g(\theta)}{\int p_\eta(x)g(\eta)d\eta}.$$

**Beispiel 2.38 (Beispiel Bern).** Wieder ist  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  durch (Bern 0) gegeben. Weiter sei die a-priori-Verteilung die Beta-Verteilung  $\pi = \beta(k, l)$ , d.h.  $\pi = p_{kl} \cdot \lambda$  mit<sup>3</sup>

$$p_{kl}(x) = \frac{\Gamma(k)\Gamma(l)}{\Gamma(k+l)} x^{k-1}(1-x)^{l-1} \sim x^{k-1}(1-x)^{l-1},$$

wobei wir mit  $\sim$  ausdrücken, dass der restliche Faktor unabhängig von  $x$  ist. Mit Bemerkung 2.37 folgt für die Dichte der a-posteriori-Verteilung  $\pi_x$

$$p_x(\theta) \sim \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \theta^{k-1} (1-\theta)^{l-1} = \theta^{k-1+\sum x_i} (1-\theta)^{l+n-1-\sum x_i}.$$

Dies ist also eine  $\beta(k + \sum x_i, l + n - \sum x_i)$ -Verteilung.

Das letzte Beispiel lässt sich auf allgemeine Exponentialfamilien verallgemeinern. (Aus Beispiel 2.24 wissen wir, dass das letzte Beispiel eine solche Verteilungsklasse behandelt.)

**Proposition 2.39 (Konjugierte Familie bei Exponentialfamilie).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  eine  $k$ -parametrische Exponentialfamilie (mit  $c, t, d, h$ ). Weiter sei  $\pi_s = p_s \cdot \lambda^k$  eine a-priori-Verteilung mit

$$p_s(\theta) = \frac{\exp\left(\sum_{j=1}^k c_j(\theta)s_j - s_{k+1}d(\theta)\right)}{\int \exp\left(\sum_{j=1}^k c_j(\eta)s_j - s_{k+1}d(\eta)\right)\lambda^k(d\eta)}$$

gegeben. Dann ist die a-posteriori-Verteilung gerade  $P(\Theta \in \cdot | X = x) = p_x \cdot \lambda^k$  mit

$$p_x(\theta) = p_{(s_1+t_1, \dots, s_k+t_k, s_{k+1}+1)}(\theta).$$

*Beweis.* Wir schreiben  $a \sim_x b$ , falls  $a/b$  nur von  $x$  abhängt (d.h.  $a$  und  $b$  sind proportional). Es gilt

$$\begin{aligned} p_x(\theta) &\sim_x p_\theta(x)\pi_s(x) \sim_x \exp\left(\sum_{j=1}^k c_j(\theta)(t_j(x) + s_j) - (s_{k+1} + 1)d(\theta)\right) \\ &\sim_x p_{(s_1+t_1, \dots, s_k+t_k, s_{k+1}+1)}(\theta). \end{aligned}$$

□

**Beispiel 2.40 (A-posteriori-Verteilung bei der Normalverteilung).** Wir betrachten das Normalverteilungsmodell bei bekanntem  $\sigma^2$ , gegeben durch (Norm 1b). Angenommen, die a-priori-Verteilung  $\pi$  für  $\theta$  ist selbst eine Normalverteilung, nämlich  $\mathcal{N}(a, b^2)$  für  $a, b \in \mathbb{R}$ . Wie

<sup>3</sup>Für die Gamma-Funktion gilt etwa  $\Gamma(k) = (k-1)!$  für  $k = 1, 2, \dots$  sowie  $\Gamma(x+1) = x\Gamma(x)$  für  $x \in \mathbb{R}$ .

sieht dann die a-posteriori-Verteilung aus? Und ist diese um den wahren Wert  $\theta$  konzentriert? In der Übung wird gezeigt werden, dass dies wieder eine Normalverteilung  $\mathcal{N}(\alpha, \beta)$  ist mit

$$\alpha = \frac{1}{\sigma^2/(nb^2) + 1} \bar{x} + \frac{\sigma^2/(nb^2)}{\sigma^2/(nb^2) + 1} a,$$

$$\beta = \frac{1}{n/\sigma^2 + 1/b^2}.$$

Insbesondere ist für große  $n$  die a-posteriori-Verteilung um  $\bar{x}$  konzentriert.

### 3 Entscheidungstheorie

Statistische Fragestellungen formuliert man oft als Entscheidungsprobleme. Diese zeichnen sich dadurch aus, dass aufgrund von Daten in einem (statistischen) Modell immer eine Entscheidung über die verwendeten Parameter getroffen werden muss. Bei einem statistischen Test ist diese Entscheidung etwa, ob der Parameter größer oder kleiner als ein vorgegebener Wert ist, bei einem Schätzproblem fällt eine Entscheidung über die (vermutete) Größe eines Parameters.

#### 3.1 Einführung

Zentral sind bei Entscheidungen die Begriffe des Entscheidungsraumes und der Entscheidungsfunktion. Um über die Qualität der Entscheidungsfunktion zu urteilen, gibt es außerdem eine Verlustfunktion.

**Definition 3.1 (Entscheidungsraum, -funktion, Verlustfunktion).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell.

1. Für einen Entscheidungsraum  $\aleph$  heißt jede messbare Abbildung  $d : E \rightarrow \aleph$  nicht-randomisierte Entscheidungsfunktion. Eine (randomisierte) Entscheidungsfunktion ist ein Markov-Kern  $\delta(\cdot, \cdot)$  von  $E$  nach  $\mathcal{B}(\aleph)$ . Die Menge der nicht-randomisierten Entscheidungsfunktionen bezeichnen wir mit  $\mathcal{D}_{nr}$ , die Menge der randomisierten Entscheidungsfunktionen mit  $\mathcal{D}$ . Für  $d \in \mathcal{D}_{nr}$  sei  $\delta_d(x, A) := \mathbb{1}_{d(x) \in A}$  die zugehörige randomisierte Entscheidungsfunktion.
2. Eine Verlustfunktion ist eine messbare Abbildung  $\ell : \mathcal{P} \times \aleph \rightarrow \mathbb{R}_+$ .
3. Ein statistisches Entscheidungsproblem ist ein Tripel  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \aleph, \ell)$  aus einem statistischen Modell, einem Entscheidungsraum und einer Verlustfunktion.
4. Für ein statistisches Entscheidungsproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \aleph, \ell)$  heißt die Abbildung

$$R_\delta(\theta) := R(\theta, \delta) := \mathbb{E}_\theta \left[ \int \ell(\theta, y) \delta(X, dy) \right] \quad (3.1)$$

Risikofunktion. Für eine nicht-randomisierte Entscheidung ist dies gleich

$$R_d(\theta) := R(\theta, d) := \mathbb{E}_\theta [\ell(\theta, d(X))].$$

Die Menge  $\mathcal{R} := \{R_\delta : \delta \in \mathcal{D}\}$  heißt Risikomenge.

**Bemerkung 3.2 (Interpretation).** In einem statistischen Modell steht  $X$  für die (zufällig entstandenen) Daten. Bei Vorliegen des (unbekannten Parameters)  $\theta \in \mathcal{P}$  sind diese nach  $\mathbb{P}_\theta$  verteilt. Für eine Entscheidungsfunktion  $\delta$  ist nun  $\delta(x, A)$  die Wahrscheinlichkeit, sich für ein  $a \in A$  zu entscheiden, wenn  $X = x$  vorliegt. (Bei einer nicht-randomisierten Entscheidungsfunktion  $d$  ist  $d(x) = a$  die Wahl der Entscheidung  $a$  im Entscheidungsraum  $\aleph$ .) Bei einer Entscheidung für  $a$  hat man, falls  $\theta$  vorliegt, einen Verlust zu verzeichnen. Dieser wird mit  $\ell(\theta, a)$  bezeichnet. Da die Daten als zufällig angesehen werden, kann man sich fragen, welchen Verlust man denn (bei Vorliegen von  $\theta$  und für die Entscheidung  $\delta$ ) erwartet. Dies ist gerade die Risikofunktion  $R_\delta(\theta)$ .

**Bemerkung 3.3 (Punkt-Schätzproblem).** 1. Aus der Vorlesung *Stochastik* bekannt ist folgendes Problem:

Ein Versuch, der entweder einen Erfolg oder einen Misserfolg bringt, wird  $n$ -mal wiederholt, wobei  $S = k$ -mal ein Erfolg eintritt. Nun soll die Wahrscheinlichkeit für einen Erfolg (in einem der  $n$  Versuchen) geschätzt werden.

Dies erinnert an das statistische Modell aus Beispiel *Bern*. Wir setzen hier  $\mathcal{P} = [0, 1]$ ,  $\mathbb{P}_\theta$  wie in (*Bern* 0),  $S = X_1 + \dots + X_n$  und  $\aleph = \mathcal{P}$ . Eine offensichtlich Wahl für eine nicht-randomisierte Entscheidungsfunktion ist  $d(x_1, \dots, x_n) = (x_1 + \dots + x_n)/n =: \bar{x}$ . Als Verlustfunktion könnte etwa  $\ell(\theta, a) = (\theta - a)^2$  dienen, also  $\ell(\theta, d(x_1, \dots, x_n)) = (\theta - \bar{x})^2$ . Die Risikofunktion berechnet sich dann zu

$$R_d(\theta) = \mathbb{E}_\theta[(\theta - \bar{X})^2] = \mathbb{V}_\theta[\bar{X}] = \frac{\theta(1 - \theta)}{n}.$$

2. Allgemeiner ist ein Punkt-Schätzproblem ein statistisches Entscheidungsproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \aleph, \ell)$  mit dem statistischen Raum  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$ . Im Allgemeinen wollen wir nicht  $\theta$  schätzen, sondern  $g(\theta)$  für eine Funktion  $g : \mathcal{P} \rightarrow \aleph$ , wobei  $\aleph$  auch der Entscheidungsraum ist. (Wie hierbei die Entscheidungsfunktion  $\delta$  (oder  $d$  im nicht-randomisierten Fall) aussieht, hängt vom jeweiligen Beispiel ab.) Verlustfunktionen, die der Bemessung einer falschen Entscheidung dienen sollen, sind (im Falle eines normierten Raumes  $\aleph$ ) etwa der

$$\text{Laplace-Verlust } \ell(\theta, a) := |a - g(\theta)|,$$

$$\text{Gauß-Verlust } \ell(\theta, a) := |a - g(\theta)|^2,$$

$$\text{0-1-Verlust } \ell(\theta, a) := 1_{|a - g(\theta)| > \varepsilon}.$$

Die Risikofunktion  $R_\delta(\theta)$  ist dann der (unter  $\mathbb{P}_\theta$ ) erwartete Verlust unter der (etwa nicht-randomisierten) Entscheidungsfunktion  $d$ , also beim Laplace-Verlust etwa

$$R_\delta(\theta) = \mathbb{E}_\theta[|d(X) - g(\theta)|].$$

Man nennt Entscheidungsfunktionen bei solchen Schätzproblemen (*Punkt-*)*Schätzer* und  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \aleph, g, \ell)$  ein Schätzproblem.

**Bemerkung 3.4 (Bereichs-Schätzproblem).** 1. Aus der Vorlesung *Stochastik* bekannt ist folgendes Problem:

Für eine Stichprobe  $x_1, \dots, x_n$  normalverteilter Daten (d.h.  $X_1, \dots, X_n$  sind unabhängig und identisch nach  $\mathcal{N}(\theta, \sigma^2)$  verteilt) bei bekannter Varianz  $\sigma^2$ . Es wird ein von den Daten abhängiger Bereich gesucht, so dass  $\theta$  mit Wahrscheinlichkeit  $1 - \alpha$  (etwa für

$\alpha = 5\%$ ) in diesem liegt. Es ist naheliegend, dass dieses Intervall symmetrisch um  $\bar{x} := (x_1 + \dots + x_n)/n$  liegt.

Dies erinnert an das statistische Modell *Norm* mit bekannter Varianz  $\sigma^2$ , d.h. an das statistische Modell  $(X, \{\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2) : \theta \in \mathbb{R}\})$ , so dass (*Norm 0*) gilt. Um das Intervall aufzustellen, verwenden wir, dass  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ , also  $(\bar{X} - \mu)/\sqrt{\sigma^2/n} \sim \mathcal{N}(0, 1)$ . Sei  $q_\alpha$  das  $\alpha$ -Quantil der Standardnormalverteilung<sup>4</sup>. Dann gilt

$$\begin{aligned} 1 - \alpha &= \mathbb{P}_\theta(q_{\alpha/2} \leq (\bar{X} - \theta)/\sqrt{\sigma^2/n} \leq q_{1-\alpha/2}) \\ &= \mathbb{P}_\theta(|\bar{X} - \theta| \leq q_{1-\alpha/2}\sqrt{\sigma^2/n}) \\ &= \mathbb{P}_\theta(\bar{X} - q_{1-\alpha/2}\sqrt{\sigma^2/n} \leq \theta \leq \bar{X} + q_{1-\alpha/2}\sqrt{\sigma^2/n}) \end{aligned}$$

Der gesuchte Datenbereich ist also  $\{y : \bar{x} - q_{1-\alpha/2}\sqrt{\sigma^2/n} \leq \bar{y} \leq \bar{x} + q_{1-\alpha/2}\sqrt{\sigma^2/n}\}$ .

2. Allgemeiner ist ein Bereichs-Schätzproblem ein statistisches Entscheidungsproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  mit dem statistischen Raum  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$ . Wieder wollen wir nicht  $\theta$  schätzen, sondern  $g(\theta)$  für eine Funktion  $g : \mathcal{P} \rightarrow \Upsilon$ , wobei hier nun  $\mathfrak{N} = \mathcal{B}(\Upsilon)$  ist. (Übrigens ist damit  $\mathfrak{N}$  kein metrischer Raum, und man muss die Messbarkeit der folgenden Abbildungen überprüfen.) Eine nicht-randomisierte Entscheidungsfunktion ist hier wieder eine Funktion  $d : E \rightarrow \mathfrak{N}$ . (Nun muss etwa für alle  $g \in \Upsilon$  die Bedingung  $\{x : g \in d(x)\} \in \mathcal{B}(E)$  gelten.) Als Verlustfunktion bietet sich die Wahl

$$\ell(\theta, B) := \begin{cases} 1, & g(\theta) \notin B, \\ 0, & g(\theta) \in B \end{cases}$$

an. Das Risiko ist somit

$$R_d(\theta) = \mathbb{E}_\theta[\ell(\theta, d(X))] = \mathbb{P}_\theta[g(\theta) \notin d(X)].$$

Entscheidungsfunktionen in einer solchen Situation heißen auch *Bereichs-Schätzer*.

**Bemerkung 3.5 (Test-Problem).** 1. Aus der Vorlesung *Stochastik* bekannt ist folgendes Problem (siehe *einfacher t-Test*):

Für eine Stichprobe  $x_1, \dots, x_n$  normalverteilter Daten (d.h.  $X_1, \dots, X_n$  sind unabhängig und identisch nach  $\mathcal{N}(\theta = (\mu, \sigma^2))$  verteilt) bei unbekannter Varianz  $\sigma^2$ . Es soll getestet werden, ob die (Null-)Hypothese  $\mu = \mu_0$  aufgrund der Daten verworfen werden kann oder nicht. Dabei soll die Wahrscheinlichkeit, die Nullhypothese irrtümlicherweise zu verwerfen, höchstens ein vorgegebenes  $\alpha$  sein.

Dies erinnert an das statistische Modell *Norm*, d.h. an das statistische Modell  $(X, \{\mathbb{P}_\theta = \mathcal{N}(\theta) : \theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+\})$ , so dass (*Norm 0*) gilt. Um die Hypothese  $H_0 : \mu = \mu_0$  gegen  $H_1 : \mu \neq \mu_0$  zu testen, verwendet man, dass für  $s^2(X) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  die Statistik

$$T := t(X) = \frac{\bar{X} - \mu}{\sqrt{s^2(X)/n}}$$

unter  $\mathbb{P}_{(\mu, \sigma^2)}$  (für jedes  $\sigma^2 \in \mathbb{R}_+$ ) nach  $t_{n-1}$ -verteilt ist.<sup>5</sup> Man stellt hier den kritischen (oder Ablehnungs-)Bereich  $C \subseteq \mathbb{R}$  (der nur von  $\alpha$  abhängt) so auf, dass  $H_0$  gerade dann

<sup>4</sup>Es ist also für  $Z \sim \mathcal{N}(0, 1)$  gerade  $\mathbb{P}(Z \leq q_\alpha) = \alpha$ . Insbesondere ist auch  $q_\alpha = -q_{1-\alpha}$ .

<sup>5</sup>Dies ist eine studentische  $t$ -Verteilung mit  $n - 1$  Freiheitsgraden. Ihr  $\alpha$ -Quantil bezeichnen wir mit  $t_{n-1, \alpha}$ .

abgelehnt wird, wenn  $t(X) \in C$  ist. Es soll außerdem  $\mathbb{P}_{(\mu_0, \sigma^2)}(t(X) \in C) \leq \alpha$  gelten. In diesem Beispiel ergibt sich mit  $C = (-\infty, t_{n-1, \alpha/2}) \cup (t_{n-1, 1-\alpha/2}, \infty)$  gerade, dass

$$\mathbb{P}_{(\mu_0, \sigma^2)}(t(X) \in C) = \mathbb{P}_{(\mu_0, \sigma^2)}(t(X) \leq t_{n-1, \alpha/2}) + \mathbb{P}_{(\mu_0, \sigma^2)}(t(X) \geq t_{n-1, 1-\alpha/2}) = \alpha,$$

was ja gefordert war.

2. Allgemeiner ist ein Test-Problem ein statistisches Entscheidungsproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  mit dem statistischen Raum  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  mit  $\mathfrak{N} = \{H_0, H_1\}$ . Die (randomisierte) Entscheidungsfunktion  $\delta$  ist eindeutig durch eine Funktion  $\varphi : E \rightarrow [0, 1]$  mittels  $\delta(x, \{H_1\}) := \varphi(x)$  und  $\delta(x, \{H_0\}) := 1 - \varphi(x)$  bestimmt. (Hier ist  $\varphi(x)$  die Wahrscheinlichkeit, sich bei Daten  $x$  für die Hypothese  $H_1$  zu entscheiden.) Die Entscheidungsfunktion  $\delta$  ist genau dann nicht-randomisiert, falls  $\varphi = 1_C$  für ein geeignetes  $C \in \mathcal{B}(E)$ . Als Verlustfunktion bietet sich die Neyman-Pearson-Verlustfunktion an. Hierbei ist für  $\mathcal{P}_0, \mathcal{P}_1$  mit  $\mathcal{P}_0 \uplus \mathcal{P}_1 = \mathcal{P}$  gerade

$$\ell(\theta, H_1) = \begin{cases} 0, & \theta \in \mathcal{P}_1 & \text{richtige Entscheidung,} \\ 1, & \theta \in \mathcal{P}_0 & \text{Fehler 1. Art: Entscheidung für } H_1, \text{ aber } H_0 \text{ trifft zu.} \end{cases}$$

$$\ell(\theta, H_0) = \begin{cases} 0, & \theta \in \mathcal{P}_0 & \text{richtige Entscheidung,} \\ 1, & \theta \in \mathcal{P}_1 & \text{Fehler 2. Art: Entscheidung für } H_0, \text{ aber } H_1 \text{ trifft zu.} \end{cases}$$

Die Hypothese  $H_i$  besteht dabei gerade daraus, dass eine der Verteilungen aus  $\mathcal{P}_i$  (auf die Daten) zutrifft,  $i = 0, 1$ . Die Risikofunktion ist damit

$$\begin{aligned} R_\delta(\theta) &= \mathbb{E}_\theta[\ell(\theta, H_0)\delta(X, \{H_0\}) + \ell(\theta, H_1)\delta(X, \{H_1\})] \\ &= \mathbb{1}_{\theta \in \mathcal{P}_1} \mathbb{E}_\theta[1 - \varphi(X)] + \mathbb{1}_{\theta \in \mathcal{P}_0} \mathbb{E}_\theta[\varphi(X)] = \begin{cases} 1 - \mathbb{E}_\theta[\varphi(X)], & \theta \in \mathcal{P}_1, \\ \mathbb{E}_\theta[\varphi(X)], & \theta \in \mathcal{P}_0. \end{cases} \end{aligned}$$

Wie wir sehen, wird das Risiko eindeutig durch die *Gütefunktion*

$$\beta_\varphi : \theta \mapsto \mathbb{E}_\theta[\varphi(X)]$$

bestimmt. Im Fall einer nicht-randomisierten Entscheidungsfunktion  $\varphi = 1_C$  ist also

$$R_\delta(\theta) = \begin{cases} \mathbb{P}_\theta(X \notin C), & \theta \in \mathcal{P}_1, \\ \mathbb{P}_\theta(X \in C), & \theta \in \mathcal{P}_0. \end{cases}$$

Bei diesem allgemeinen Vorgehen fällt auf, dass nirgends das Signifikanzniveau  $\alpha$  aus obigem Beispiel eingeht. Das liegt daran, dass wir zunächst alle Entscheidungsfunktionen zugelassen haben. Sinnvoller ist es jedoch, sich auf eine Klasse von Entscheidungsfunktionen einzuschränken, etwa auf solche  $\delta$ , für die  $R_\delta(\theta) \leq \alpha$  für  $\theta \in \mathcal{P}_0$  gilt.

Man nennt Entscheidungsfunktionen  $\delta$  (oder  $\varphi$ ) bei solchen Testproblemen auch *Tests*.

**Bemerkung 3.6 (Randomisierter Test).** Wir geben noch das Beispiel eines randomisierten Tests an. Sei hierzu  $(X, \{\mathbb{P}_\theta : \theta \in [0, 1]\})$  das statistische Modell aus Beispiel 1.6 sowie  $H_0 = \{\theta \leq 0.5\}$  und  $H_1 = \{\theta > 0.5\}$ . Für ein  $\alpha \in (0, 1)$  wollen wir einen Test  $\varphi$  mit  $\mathbb{E}_\theta[\varphi(X)] \leq \alpha$  für  $\theta \in H_0$  angeben, d.h. die Wahrscheinlichkeit, sich für die Alternative zu entscheiden, soll bei Vorliegen von  $\theta \in H_0$  höchstens  $\alpha$  betragen. Man sagt auch, der Fehler 1.



Art soll höchstens  $\alpha$  sein. Sei hierzu  $q_{\theta,1-\alpha}$  ein  $1 - \alpha$ -Quantil einer  $B(n, \theta)$ -Verteilung<sup>6</sup>. Dann könnten wir etwa

$$\varphi(x) = \begin{cases} 0, & \sum_{i=1}^n x_i \leq q_{0.5,1-\alpha}, \\ 1, & \sum_{i=1}^n x_i > q_{0.5,1-\alpha} \end{cases}$$

setzen. Für  $\theta \in H_0$  gilt dann

$$\mathbb{E}_\theta[\varphi(X)] = \mathbb{P}_\theta \left[ \sum_{i=1}^n X_i > q_{0.5,1-\alpha} \right] \leq \mathbb{P}_{0.5} \left[ \sum_{i=1}^n X_i > q_{0.5,1-\alpha} \right] \leq \alpha.$$

Schön wäre es, wenn wir sogar ein  $\psi \geq \varphi$  angeben können, für das ebenfalls  $\mathbb{E}_\theta[\psi(X)] \leq \alpha$  für  $\theta \in H_0$  gilt. Dies können wir erreichen, indem wir

$$\psi(x) = \begin{cases} 0, & \sum_{i=1}^n x_i < q_{0.5,1-\alpha}, \\ p := \frac{\mathbb{P}_{0.5} \left[ \sum_{i=1}^n x_i \leq q_{0.5,1-\alpha} \right] - (1-\alpha)}{\mathbb{P}_{0.5} \left[ \sum_{i=1}^n x_i = q_{0.5,1-\alpha} \right]}, & \sum_{i=1}^n x_i = q_{0.5,1-\alpha}, \\ 1, & \sum_{i=1}^n x_i > q_{0.5,1-\alpha}. \end{cases}$$

Das bedeutet: falls  $\sum_{i=1}^n x_i = q_{0.5,1-\alpha}$ , so entscheiden wir uns mit Wahrscheinlichkeit  $p$  für die Alternative, und mit  $1 - p$  für die Nullhypothese. Die Entscheidung ist also zufällig oder randomisiert. Dann ist für  $\theta \in H_0$

$$\begin{aligned} \mathbb{E}_\theta[\psi(X)] &\leq \mathbb{E}_{0.5}[\psi(X)] = 1 - \mathbb{P}_{0.5} \left[ \sum_{i=1}^n X_i \leq q_{0.5,1-\alpha} \right] + p \mathbb{P}_{0.5} \left[ \sum_{i=1}^n X_i = q_{0.5,1-\alpha} \right] \\ &= 1 - \mathbb{P}_{0.5} \left[ \sum_{i=1}^n X_i \leq q_{0.5,1-\alpha} \right] + \mathbb{P}_{0.5} \left[ \sum_{i=1}^n X_i \leq q_{0.5,1-\alpha} \right] - (1 - \alpha) = \alpha, \end{aligned}$$

also ist der Fehler 1. Art ebenfalls höchstens  $\alpha$ .

Es sollte klar sein, dass *gute* Entscheidungskriterien gerade solche sind, die eine kleine Risikofunktion aufweisen. Allerdings gibt es verschiedene Möglichkeiten dafür, dass eine Entscheidungsfunktion *klein* ist. Mit diesen werden wir uns in Abschnitt 3.3 beschäftigen. Wir geben nur noch ein Beispiel, in dem klar wird, dass naiv angegebene Entscheidungsfunktionen auch *schlecht* sein können.

### Beispiel 3.7 (Schätzung des Lageparameters einer Uniformverteilung).

Sei  $U(\theta - 1/2, \theta + 1/2)$  die Uniformverteilung auf  $(\theta - 1/2, \theta + 1/2)$  und  $\mathbb{P}_\theta = U(\theta - 1/2, \theta + 1/2)^n$  das  $n$ -fache Produkt,  $\theta \in \mathbb{R}$ . Wir wollen für das statistische Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathbb{R}\})$  den Lageparameter  $\theta$  schätzen. Hierzu geben wir zwei verschiedene, nicht-randomisierte Entscheidungsfunktionen an. Wir wählen (mit  $g(\theta) := \theta$ )<sup>7</sup>

$$d_1(X) = \bar{X}, \quad d_2(X) = \frac{1}{2}(X_{(1)} + X_{(n)})$$

<sup>6</sup>Wir erinnern daran, dass für eine Verteilung  $\mathbb{P}$  mit Verteilungsfunktion  $F$  das  $\alpha$ -Quantil durch jede Zahl  $x$  mit  $F(x-) \leq \alpha \leq F(x)$  gegeben ist.

<sup>7</sup>Für einen Vektor  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  ist  $x_{(i)}$  das  $i$ t-kleinste Element. Insbesondere ist also  $x_{(1)}$  das kleinste und  $x_{(n)}$  das größte Element. Die Zahlen  $x_{(1)}, \dots, x_{(n)}$  heißen auch *Ordnungsstatistiken*.

und verwenden den Gauß-Verlust. Das Risiko berechnet sich für  $d_1$  zu

$$R_{d_1}(\theta) = \mathbb{E}_\theta[(\bar{X} - \theta)^2] = \mathbb{V}_\theta[\bar{X}] = \frac{1}{n} \mathbb{V}_\theta[X_1] = \frac{1}{12n}.$$

Für  $d_2$  benötigen wir ein paar Vorüberlegungen. Um die gemeinsame Verteilung von  $X_{(1)}$  und  $X_{(n)}$  zu berechnen, schreiben wir (für  $\theta = 1/2$ )

$$\mathbb{P}_{1/2}(X_{(1)} > x, X_{(n)} < y) = (y - x)^n.$$

Also hat  $(X_{(1)}, X_{(n)})$  die gemeinsame Dichte  $(x, y) \mapsto n(n-1)(y-x)^{n-2}$ . Daraus berechnen wir

$$\begin{aligned} \mathbb{E}_{1/2}[X_{(1)}] &= 1 - \mathbb{E}_{1/2}[1 - X_{(1)}] = 1 - \int_0^1 n(1-x)^n dx = \frac{1}{n+1}, \\ \mathbb{E}_{1/2}[X_{(n)}] &= \int_0^1 nyy^{n-1} dy = \frac{n}{n+1}, \\ \mathbb{V}_{1/2}[X_{(1)}] &= \mathbb{V}_{1/2}[1 - X_{(1)}] = \int_0^1 n(1-x)^{n+1} dx - \frac{n^2}{(n+1)^2} = \frac{n}{n+2} - \frac{n^2}{(n+1)^2} \\ &= \frac{n}{(n+2)(n+1)^2}, \\ \mathbb{V}_{1/2}[X_{(n)}] &= \frac{n}{(n+2)(n+1)^2}, \\ \text{COV}_{1/2}[X_{(1)}, X_{(n)}] &= \int_0^1 \int_0^y n(n-1)xy(y-x)^{n-2} dx dy - \frac{n}{(n+1)^2} \\ &= \int_0^1 \int_0^y ny(y-x)^{n-1} dx dy - \frac{n}{(n+1)^2} \\ &= \int_0^1 y^{n+1} dy - \frac{n}{(n+1)^2} = \frac{1}{n+2} - \frac{1}{(n+1)^2} = \frac{1}{(n+1)^2(n+2)}, \end{aligned}$$

wobei die Varianzen und Kovarianz unabhängig von  $\theta$  sind. Wir erhalten

$$\begin{aligned} R_{d_1}(\theta) &= \mathbb{E}_\theta\left[\left(\frac{1}{2}(X_{(1)} + X_{(n)}) - \theta\right)^2\right] = \mathbb{V}_\theta\left[\frac{1}{2}(X_{(1)} + X_{(n)})\right] \\ &= \frac{1}{4}(\mathbb{V}[X_{(1)}] + \mathbb{V}[X_{(n)}] + 2\text{COV}[X_{(1)}, X_{(n)}]) \\ &= \frac{1}{2}\left(\frac{n}{(n+2)(n+1)^2} + \frac{1}{(n+1)^2(n+2)}\right) = \frac{1}{2(n+1)(n+2)}. \end{aligned}$$

Wir sehen also, dass das Risiko von  $d_2$  gerade für große  $n$  deutlich kleiner ist als das von  $d_1$ .

### 3.2 Die Rolle suffizienter Statistiken

Suffiziente Statistiken enthalten alle wichtigen Informationen über die Daten. Deswegen ist es sinnvoll, dass Entscheidungsfunktionen nur von solchen Statistiken abhängen.

**Proposition 3.8.** *Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem,  $\delta$  eine (randomisierte) Entscheidungsfunktion (also ein Markov-Kern von  $E$  nach  $\mathcal{B}(\mathfrak{N})$ ) und  $R_\delta$  ihre Risikofunktion. Ist  $T = t(X)$  (für  $t : E \rightarrow E'$ ) eine suffiziente Statistik, so gibt es einen Markov-Kern  $\varepsilon$  von  $E'$  nach  $\mathcal{B}(\mathfrak{N})$ , so dass  $\varepsilon \circ t$  (mit  $\varepsilon \circ t(x, A) := \varepsilon(t(x), A)$ ) eine Entscheidungsfunktion ist mit  $R_{\varepsilon \circ t} = R_\delta$ .*

*Beweis.* Für  $A \in \mathcal{B}(\mathfrak{N})$  setzen wir

$$\varepsilon(t, A) := \mathbb{E}_\theta[\delta(X, A)|T = t].$$

(Da  $T$  suffizient ist, hängt die rechte Seite nicht von  $\theta$  ab. Wir unterdrücken den Subskript  $\theta$  im Folgenden.) Deshalb gilt für integrierbare Funktionen  $h : \mathfrak{N} \rightarrow \mathbb{R}$ , dass

$$\mathbb{E} \left[ \int h(a) \delta(X, da) \middle| T = t \right] = \int h(a) \varepsilon(t, da).$$

(Zunächst überprüft man die Aussage mit Indikatorfunktionen, dann mit einfachen Funktionen und anschließend mit allgemeinen messbaren Funktionen  $h$ .) Nach Definition gilt dann

$$\begin{aligned} R_{\varepsilon \circ t}(\theta) &= \mathbb{E}_\theta \left[ \int \ell(\theta, a) \varepsilon(t(X), da) \right] = \mathbb{E}_\theta \left[ \mathbb{E} \left[ \int \ell(\theta, a) \delta(X, da) \middle| T \right] \right] \\ &= \mathbb{E}_\theta \left[ \int \ell(\theta, a) \delta(X, da) \right] = R_\delta(\theta). \end{aligned}$$

□

Man beachte, dass auch bei nicht-randomisierten  $\delta$  der Markov-Kern  $\varepsilon$  (und damit  $\varepsilon \circ t$ ) randomisiert sein kann.

**Proposition 3.9.** *Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem. Sei  $\mathfrak{N} \subseteq \mathbb{R}^m$  konvex und  $\ell(\theta, \cdot)$  für alle  $\theta \in \mathcal{P}$  eine konvexe Verlustfunktion. Für eine (randomisierte) Entscheidungsfunktion  $\delta$  definieren wir die nicht-randomisierte Entscheidungsfunktion*

$$d(x) := \int a \delta(x, da),$$

zumindest für  $x \in E$ , für die dieses Integral existiert. Für solche  $x \in E$  und alle  $\theta \in \mathcal{P}$  ist dann

$$\ell(\theta, d(x)) \leq \int \ell(\theta, a) \delta(x, da).$$

*Beweis.* Zunächst ist  $d(x) \in \mathfrak{N}$ , da  $\mathfrak{N}$  konvex ist. Weiter gilt mit der Jensen'schen Ungleichung

$$\ell(\theta, d(x)) = \ell\left(\theta, \int a \delta(x, da)\right) \leq \int \ell(\theta, a) \delta(x, da).$$

□

Man kann mit Hilfe suffizienter Statistiken Entscheidungsfunktionen besser machen. Vor allem für Schätzprobleme wird folgendes Resultat schöne Anwendungen haben.

**Theorem 3.10 (Rao-Blackwell).** *Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem. Sei  $\mathfrak{N} \subseteq \mathbb{R}^m$  konvex und  $\ell(\theta, \cdot)$  für alle  $\theta \in \mathcal{P}$  eine konvexe Verlustfunktion,  $T = t(X)$  eine suffiziente Statistik und  $d$  eine nicht-randomisierte Entscheidungsfunktion mit  $\mathbb{E}_\theta[|d(X)|] < \infty$  für alle  $\theta \in \mathcal{P}$  sowie  $R_d$  ihre Risikofunktion. Für*

$$e(t) := \mathbb{E}_\theta[d(X)|T = t]$$

(wobei die rechte Seite unabhängig von  $\theta$  ist) gilt dann für alle  $\theta \in \mathcal{P}$

$$R_{e \circ t}(\theta) \leq R_d(\theta).$$

*Beweis.* Wir setzen

$$\varepsilon(t, A) := \mathbb{P}(d(X) \in A | T = t),$$

also

$$\int h(a)\varepsilon(t, da) = \mathbb{E}[h(d(X)) | T = t],$$

falls das Integral existiert. (Wieder zeigt man das, indem man zuerst einfache Funktionen einsetzt und danach durch ein Approximationsargument messbare Funktionen.) Daraus folgt insbesondere

$$\int a\varepsilon(t, da) = \mathbb{E}[d(X) | T = t] = e(t).$$

Dann gilt wegen Proposition 3.9

$$\ell(\theta, e \circ t(x)) \leq \int \ell(\theta, a)\varepsilon(t(x), da),$$

also

$$R_{e \circ t}(\theta) \leq R_\varepsilon(\theta) = \mathbb{E}_\theta \left[ \int \ell(\theta, a)\varepsilon(T, da) \right] = \mathbb{E}_\theta[\mathbb{E}_\theta[\ell(\theta, d(X)) | T]] = R_d(\theta).$$

□

**Beispiel 3.11 (Beispiel Unif).** Sei  $(X, \{\mathcal{P}_\theta : \theta \in (0, 1)\})$  wie in Beispiel 1.8,  $\mathfrak{N} = (0, 1)$ ,  $g(\theta) = \theta$  und  $\ell$  der Gauß-Verlust  $\ell(\theta, a) = (\theta - a)^2$ . Wir betrachten den Schätzer  $d(x) = 2\bar{x}$ . Für diesen ist immerhin  $\mathbb{E}_\theta[d(X)] = \theta$ . Wir wissen aus Beispiel 2.8, dass  $T = t(X) = \max_{i=1, \dots, n} X_i$  suffizient ist. Wir verwenden Theorem 3.10 und schreiben aus Symmetriegründen

$$\mathbb{E}_\theta[2\bar{X} | t(X)] = 2\mathbb{E}_\theta[X_1 | t(X)] = 2\left(\frac{1}{n}t(X) + \frac{n-1}{n}\frac{t(X)}{2}\right) = \frac{n+1}{n}t(X).$$

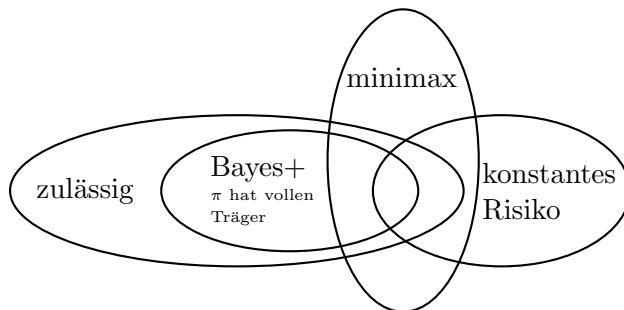
Nun lesen wir ab, dass der Schätzer

$$d'(x) := \frac{n+1}{n} \max_{i=1, \dots, n} x_i$$

eine kleinere Risikofunktion hat.

### 3.3 Zulässige, Bayes, Minimax-Entscheidungsfunktionen

In diesem Abschnitt geht es um den Vergleich von Entscheidungsfunktionen mittels deren Risikofunktionen, sowie um bestimmte Optimalitätskriterien. Wie werden hier speziell zulässige, minimax und Bayes-Entscheidungsfunktionen betrachten. Oftmals macht es dabei keinen Sinn, optimale Entscheidungskriterien unter allen möglichen Entscheidungsfunktionen zu suchen, sondern nur aus einer Teilmenge. (Wir erinnern daran, dass wir die Menge aller Entscheidungsfunktionen mit  $\mathcal{D}$  bezeichnet hatten.) Wir illustrieren einige Ergebnisse dieses Abschnitts in einer Grafik.



In Lemma 3.14 zeigen wir etwa, dass zulässige Entscheidungsfunktionen mit konstantem Risiko minimax sind. Lemma 3.19 besagt, dass Bayes-Entscheidungsfunktionen (unter einer Voraussetzung an die a-priori-Verteilung) zulässig sind. Die minimax-Eigenschaft von (allgemeinen) Bayes-Entscheidungsfunktionen wird dann mit Theorem 3.20 geklärt.

**Beispiel 3.12.** 1. Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathbb{R}\}), \mathfrak{N}, \ell)$  mit  $\mathbb{P}_\theta = \mathcal{N}(\theta, 1)$ ,  $\mathfrak{N} = \mathbb{R}$  und  $\ell$  der Gauß-Verlust. Wir wollen also den Mittelwert einer Normalverteilung (bei bekannter Varianz) schätzen, wenn wir nur einmal aus ihr ziehen. Die offensichtliche Wahl  $\delta_1(x, a) = x$  (d.h. wir schätzen den Parameter  $\theta$  durch die Beobachtung  $x$ ) führt zur Risikofunktion

$$\theta \mapsto R_{\delta_1}(\theta) = \mathbb{E}_\theta[(X - \theta)^2] = 1.$$

Hingegen führt die Wahl  $\delta_2(x, a) = b$  für ein festes  $b \in \mathbb{R}$  (d.h. wir schätzen den Parameter  $\theta$  immer durch dieselbe Zahl  $b$ , unabhängig von unserer Beobachtung) zur Risikofunktion

$$\theta \mapsto R_{\delta_2}(\theta) = \mathbb{E}_\theta[(b - \theta)^2] = (b - \theta)^2.$$

Insbesondere sehen wir, dass weder  $R_{\delta_1} \leq R_{\delta_2}$  noch  $R_{\delta_2} \leq R_{\delta_1}$  gilt. Man kann also nicht sagen, dass  $\delta_1$  besser wäre als  $\delta_2$ . Ein Ausweg hier ist es, sich nur auf erwartungstreue Schätzer einzuschränken, d.h. auf Entscheidungsfunktionen  $\delta$  mit  $\mathbb{E}_\theta[\int \delta(X, da)] = \theta$ , und unter diesen die kleinste Risikofunktion zu suchen.

2. Ganz ähnlich verhält es sich bei einem Test. Verwenden wir hierzu die Situation aus Bemerkung 3.5.2. Es ist verlockend, einfach  $C = E$  zu wählen (d.h. die Hypothese wird immer verworfen), weil dadurch zumindest für  $\theta \in \mathcal{P}_0$  die Risikofunktion verschwindet. Allerdings wird durch diese Wahl die Risikofunktion für  $\theta \in \mathcal{P}_1$  maximal. Wählt man andererseits  $C = \emptyset$ , so ist die Situation genau andersherum. Deshalb fragt man hier eher nach dem Minimum aller Risikofunktionen, die  $R_\delta(\theta) \leq \alpha$  für alle  $\theta \in \mathcal{P}_0$  erfüllen. Das bedeutet, dass man das Signifikanzniveau  $\alpha$  festlegt und damit den maximalen Fehler erster Art.

**Definition 3.13 (Zulässige, Minimax-Entscheidungsfunktionen).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem,  $\delta, \delta' \in \mathcal{D}$  Entscheidungsfunktionen und  $R_\delta, R_{\delta'}$  ihre Risikofunktionen. Weiter sei  $\mathcal{D}_0 \subseteq \mathcal{D}$ .

1. Gilt für  $\delta, \delta'$ , dass  $R_\delta \leq R_{\delta'}$  und  $R_\delta(\theta) < R_{\delta'}(\theta)$  für mindestens ein  $\theta \in \mathcal{P}$ , so sagen wir,  $\delta$  dominiert  $\delta'$ .
2. Wir sagen, die Entscheidungsfunktion  $\delta$  hat konstantes Risiko, falls  $\theta \mapsto R_\delta(\theta)$  konstant ist.

3. Gibt es für eine Entscheidungsfunktion  $\delta \in \mathcal{D}_0$  ein  $\delta' \in \mathcal{D}_0$ , das  $\delta$  dominiert, so heißt  $\delta$  auch unzulässig für  $\mathcal{D}_0$ , andersfalls zulässig für  $\mathcal{D}_0$ .
4. Die Entscheidungsfunktion  $\delta \in \mathcal{D}_0$  heißt minimax für  $\mathcal{D}_0$ , falls

$$\sup_{\theta \in \mathcal{P}} R_\delta(\theta) = \inf_{\delta' \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta).$$

Für  $\mathcal{D}$  zulässige Entscheidungsfunktionen heißen auch zulässig, für  $\mathcal{D}$  unzulässige auch unzulässig, und für  $\mathcal{D}$  minimax-Entscheidungsfunktionen heißen auch minimax-Entscheidungsfunktionen.

**Lemma 3.14 (Zusammenhänge zulässige, minimax-Entscheidungsfunktionen).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem,  $\delta \in \mathcal{D}_0 \subseteq \mathcal{D}$  eine Entscheidungsfunktion und  $R_\delta$  ihre Risikofunktion.

1. Ist  $\delta$  die einzige minimax-Entscheidungsfunktion für  $\mathcal{D}_0$ , so ist  $\delta$  zulässig für  $\mathcal{D}_0$ .
2. Ist  $\delta$  zulässig für  $\mathcal{D}_0$  mit konstantem Risiko, so ist  $\delta$  minimax für  $\mathcal{D}_0$ .

*Beweis.* 1. Angenommen,  $\delta$  ist nicht zulässig für  $\mathcal{D}_0$ . Dann gibt es ein  $\delta' \in \mathcal{D}_0$  mit ( $\delta' \neq \delta$  und)  $R_{\delta'} \leq R_\delta$  und ein  $\theta \in \mathcal{P}$  mit  $R_{\delta'}(\theta) < R_\delta(\theta)$ . Deshalb ist

$$\inf_{\varepsilon \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_\varepsilon(\theta) = \sup_{\theta \in \mathcal{P}} R_\delta(\theta) \geq \sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta) \geq \inf_{\varepsilon \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_\varepsilon(\theta).$$

Damit ist  $\delta'$  ebenfalls minimax für  $\mathcal{D}_0$  im Widerspruch zur Voraussetzung.

2. Angenommen, es ist  $R_\delta(\theta) =: c$  unabhängig von  $\theta$  und  $\delta$  ist nicht minimax für  $\mathcal{D}_0$ . Dann gibt es  $\delta' \in \mathcal{D}_0$  mit

$$\sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta) < \sup_{\theta \in \mathcal{P}} R_\delta(\theta) = c.$$

Daraus folgt  $R_{\delta'} < R_\delta$  im Widerspruch zur Zulässigkeit für  $\mathcal{D}_0$  von  $\delta$ . □

**Definition 3.15 (Bayes'sches Risiko).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem,  $\delta$  eine Entscheidungsfunktion und  $R_\delta$  ihre Risikofunktion. Sei weiter  $\pi$  ein Wahrscheinlichkeitsmaß auf  $\mathcal{P}$  und  $\Theta \sim \pi$ . Dann ist das Bayes-Risiko (bezüglich  $\pi$ ) gegeben als

$$r_\delta(\pi) := \mathbb{E}_\pi[R_\delta(\Theta)] := \int R_\delta(\theta) \pi(d\theta).$$

Für  $\mathcal{D}_0 \subseteq \mathcal{D}$  heißt  $\delta \in \mathcal{D}_0$  eine  $\mathcal{D}_0$ -Bayes-Entscheidungsfunktion bezüglich  $\pi$ , falls

$$r_\delta(\pi) \leq r_{\delta'}(\pi)$$

für alle  $\delta' \in \mathcal{D}_0$ . Eine  $\mathcal{D}$ -Bayes-Entscheidungsfunktion bezüglich  $\pi$  heißt auch Bayes-Entscheidungsfunktion bezüglich  $\pi$ .

Oftmals sind Bayes-Entscheidungsfunktionen gar nicht so schwer zu finden. Hierbei besonders hilfreich ist das nächste Resultat.

**Theorem 3.16 (Konstruktion von Bayes-Entscheidungsfunktionen).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem. Sei  $\pi$  eine a-priori-Verteilung,  $\pi_x$  die zugehörige a-posteriori-Verteilung und  $\Theta_x \sim \pi_x$ .

1. Sei  $\delta$  eine Entscheidungsfunktion, so dass für alle  $x$

$$\mathbb{E} \left[ \int \ell(\Theta_x, a) \delta(x, da) \right] = \inf_{m \in \mathcal{M}_1(\mathbb{N})} \mathbb{E} \left[ \int \ell(\Theta, a) m(da) \right].$$

Dann ist  $\delta$  eine Bayes-Entscheidungsfunktion bezüglich  $\pi$ .

2. Sei  $d \in \mathcal{D}_{nr}$  eine nicht-randomisierte Entscheidungsfunktion, so dass für alle  $x$

$$\mathbb{E}[\ell(\Theta_x, d(x))] = \inf_{a \in \mathbb{N}} \mathbb{E}[\ell(\Theta, a)].$$

Dann ist  $\delta$  eine Bayes-Entscheidungsfunktion bezüglich  $\pi$  für  $\mathcal{D}_{nr}$ .

*Beweis.* Wir zeigen nur 1. da der Beweis für 2. analog funktioniert. Wir schreiben das Bayes-Risiko für  $\Theta \sim \pi$  als

$$r_\delta(\pi) = \mathbb{E}[R_\delta(\Theta)] = \mathbb{E} \left[ \mathbb{E} \left[ \int \ell(\Theta, a) \delta(X, da) \middle| X \right] \right].$$

Da die Verteilung von  $\Theta$  gegeben  $X$  gerade  $\pi_X$  ist, die Erwartung sicher dann minimiert, wenn

$$\mathbb{E} \left[ \int \ell(\Theta, a) \delta(X, da) \middle| X = x \right] = \mathbb{E} \left[ \int \ell(\Theta_x, a) \delta(x, da) \right]$$

minimal ist. Daraus folgt die Behauptung.  $\square$

**Korollar 3.17 (Bayes-Schätzer).** Wir betrachten das Schätzproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathbb{N} = \mathbb{R}, g, \ell)$  für den Gauß-Verlust  $\ell$ , ein reguläres statistisches Modell  $\{\mathbb{P}_\theta : \theta \in \mathcal{P}\}$  mit  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  und  $\mathcal{P} \subseteq \mathbb{R}^k$ . Weiter sei  $\pi = p \cdot \lambda^k \in \mathcal{M}_1(\mathcal{P})$  die a-priori-Verteilung und  $\pi_x$  die a-posteriori-Verteilung. (Nach Bemerkung 2.37 hat diese die Dichte

$$p_x(\theta) = \frac{p_\theta(x)p(\theta)}{\int p_\eta(x)p(\eta)d\eta}.)$$

Dann ist der nicht-randomisierte Schätzer (falls das Integral existiert)

$$d(x) := \int g(\theta) \pi_x(d\theta).$$

ein Bayes-Schätzer.

*Beweis.* Nach Proposition 3.9 müssen wir nur zeigen, dass  $d$  ein Bayes-Schätzer bezüglich  $\pi$  für  $\mathcal{D}_{nr}$  ist. Nach Theorem 3.16 ist ein Bayes-Schätzer gegeben, wenn

$$\mathbb{E}[\ell(\Theta_x, d(x))] = \min_a \mathbb{E}[\ell(\Theta_x, a)].$$

Nun ist

$$\mathbb{E}[\ell(\Theta_x, a)] = \mathbb{E}[(g(\Theta_x) - a)^2] \geq \mathbb{E}[(g(\Theta_x) - \mathbb{E}[g(\Theta_x)])^2]$$

mit '=' genau dann, wenn

$$a = \mathbb{E}[g(\Theta_x)] = \int g(\theta) \pi_x(d\theta) = d(x).$$

Daraus folgt die Behauptung.  $\square$

**Beispiel 3.18 (Beispiel Norm mit  $\sigma^2$  bekannt).** Wir betrachten für (bekanntes)  $\sigma^2 > 0$  wie in (Norm 1b) das Schätzproblem  $((X, \{\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2)^n : \theta \in \mathbb{R}\}), \mathfrak{N} = \mathbb{R}, g = \text{id}, \ell)$  für den Gauß-Verlust  $\ell$ . Wir zeigen, dass für die a-priori-Verteilung  $\pi_b = \mathcal{N}(0, b^2)$  der Bayes-Schätzer durch

$$d_b(x) := \frac{nb^2}{nb^2 + \sigma^2} \bar{x}$$

gegeben ist. (Diesen liest man am besten als Konvexkombination aus dem a-priori-Schätzer 0 – gegeben durch die a-priori-Verteilung, die Erwartungswert 0 hat – und dem Mittelwert.) Aus Beispiel 2.40 können wir die a-posteriori-Verteilung ablesen. Diese ist gerade

$$\pi_x := \mathcal{N}\left(\frac{nb^2}{nb^2 + \sigma^2} \bar{x}, \frac{\sigma^2 b^2}{nb^2 + \sigma^2}\right).$$

Nach Proposition 3.17 ist nun der Bayes-Schätzer bezüglich  $\pi$  gegeben durch

$$d(x) = \int \theta \pi_x(d\theta) = \frac{nb^2}{nb^2 + \sigma^2} \bar{x}.$$

**Lemma 3.19 (Zusammenhänge zulässige, Bayes-Entscheidungsfunktionen).** *Beachte dieselbe Situation wie in Lemma 3.14. Weiter sei  $\pi \in \mathcal{M}_1(\mathcal{P})$ .*

1. Sei  $\theta \mapsto R_{\delta'}(\theta)$  für alle  $\delta \in \mathcal{D}_0$  stetig und  $\pi$  habe vollen Träger. Ist  $\delta$  Bayes-Entscheidungsfunktion bezüglich  $\pi$  für  $\mathcal{D}_0$ , so ist  $\delta$  zulässig bezüglich  $\mathcal{D}_0$ .
2. Ist  $\delta$  die einzige Bayes-Entscheidungsfunktion bezüglich  $\pi$  und  $\mathcal{D}_0$ , so ist  $\delta$  zulässig bezüglich  $\mathcal{D}_0$ .

*Beweis.* 1. Angenommen,  $\delta$  ist nicht zulässig  $\mathcal{D}_0$ . Dann gibt es ein  $\delta' \in \mathcal{D}_0$  mit ( $\delta' \neq \delta$  und)  $R_{\delta'} \leq R_\delta$  und ein  $\theta \in \mathcal{P}$  mit  $R_{\delta'}(\theta) < R_\delta(\theta)$ . Wegen der Stetigkeit von  $R_{\delta'}$  gibt es ein  $r > 0$ , so dass für  $\theta' \in B_r(\theta)$  gerade  $R_{\delta'}(\theta') < R_\delta(\theta) - r$ . Damit gilt

$$\begin{aligned} r_{\delta'}(\pi) &= \mathbb{E}_\pi[R_{\delta'}(\Theta), \Theta \in B_r(\theta)] + \mathbb{E}_\pi[R_{\delta'}(\Theta), \Theta \notin B_r(\theta)] \\ &\leq \mathbb{E}_\pi[R_\delta(\Theta)] - r\mathbb{P}(\Theta \in B_r(\theta)) < \mathbb{E}_\pi[R_\delta(\Theta)] = r_\delta(\pi) \end{aligned}$$

im Widerspruch dazu, dass  $\delta$  Bayes-Entscheidungsfunktion bezüglich  $\pi$  für  $\mathcal{D}_0$  ist.

2. Sei  $\delta'$  so, dass  $R_{\delta'} \leq R_\delta$ . Dann gilt

$$r_{\delta'}(\pi) = \mathbb{E}[R_{\delta'}(\Theta)] \leq \mathbb{E}[R_\delta(\Theta)] = r_\delta(\pi) = \inf_{\varepsilon \in \mathcal{D}_0} r_\varepsilon(\pi).$$

Damit ist auch  $\delta'$  eine Bayes-Entscheidungsfunktion für  $\mathcal{D}_0$ . Diese ist aber nach Voraussetzung eindeutig und damit  $\delta' = \delta$ . Damit ist  $\delta$  zulässig für  $\mathcal{D}_0$ .  $\square$

**Theorem 3.20 (Hodges-Lehmann).** *Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem,  $\delta, \delta_1, \delta_2, \dots \in \mathcal{D}_0 \subseteq \mathcal{D}$  Entscheidungsfunktionen,  $R_\delta, R_{\delta_1}, R_{\delta_2}, \dots$  ihre Risikofunktionen und  $\pi, \pi_1, \pi_2, \dots$  Wahrscheinlichkeitsmaße auf  $\mathcal{P}$  sowie  $\Theta \sim \pi, \Theta_1 \sim \pi_1, \dots$*

1. Ist  $\delta_n$  eine Bayes-Entscheidungsfunktion bezüglich  $\pi_n$  für  $\mathcal{D}_0$ ,  $n = 1, 2, \dots$  und

$$\sup_{\theta \in \mathcal{P}} R_\delta(\theta) \leq \limsup_{n \rightarrow \infty} r_{\delta_n}(\pi_n),$$

dann ist  $\delta$  auch eine minimax-Entscheidungsfunktion für  $\mathcal{D}_0$ .



2. Ist  $\delta$  eine Bayes-Entscheidungsfunktion bezüglich  $\pi$  für  $\mathcal{D}_0$  mit konstantem Risiko, so ist  $\delta$  auch minimax für  $\mathcal{D}_0$ .

*Beweis.* 1. Für  $\delta \in \mathcal{D}_0$  und  $k = 1, 2, \dots$  gilt

$$\sup_{\theta \in \mathcal{P}} R_\delta(\theta) \geq \mathbb{E}_{\pi_k}[R_\delta(\Theta_k)] = r_\delta(\pi_k) \geq r_{\delta_k}(\pi_k).$$

Daraus folgt

$$\inf_{\delta' \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta) \limsup_{k \rightarrow \infty} r_{\delta_k}(\pi_k) \geq \sup_{\theta \in \mathcal{P}} R_\delta(\theta) \geq \inf_{\delta' \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta),$$

woraus insbesondere  $\sup_{\theta \in \mathcal{P}} R_\delta(\theta) = \inf_{\delta' \in \mathcal{D}_0} \sup_{\theta \in \mathcal{P}} R_{\delta'}(\theta)$  folgt. Deshalb ist  $\delta$  eine minimax-Entscheidungsfunktion für  $\mathcal{D}_0$ .

2. ist ein Spezialfall von 1. Für  $\delta$  mit konstantem Risiko ist nämlich  $\sup_{\theta \in \mathcal{P}} R_\delta(\theta) = \mathbb{E}_\pi[R_\delta(\Theta)]$ . Damit ist  $\delta$  auch eine Bayes-Entscheidungsfunktion für  $\mathcal{D}_0$ .  $\square$

Als Beispiel zeigen wir nun die Optimalität (im Sinne der Zulässigkeit) des arithmetischen Mittels zur Schätzung des Erwartungswertes bei normalverteilten Daten.

**Proposition 3.21 (Zulässigkeit des arithmetischen Mittels).** *Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathbb{R}\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem mit  $\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2)^n$  für  $\sigma^2 > 0$ ,  $n \in \mathbb{N}$ , sowie  $\mathfrak{N} = \mathbb{R}$  und  $\ell$  der Gauß-Verlust für die Funktion  $g(\theta) = \theta$ . Dann ist die Entscheidungsfunktion  $d(x) := \bar{x}$  zulässig und minimax.*

*Beweis.* Die Risikofunktion von  $d$  ist

$$R_d(\theta) = \mathbb{E}_\theta[(\bar{X} - \theta)^2] = \frac{\sigma^2}{n}.$$

Angenommen,  $d$  ist nicht zulässig. Dann gibt es einen – nach Proposition 3.9 nicht-randomisierte – Schätzer  $d'$ , der  $d$  dominiert. Es gilt also  $R_{d'} \leq \sigma^2/n$  und es gibt ein  $\theta \in \mathbb{R}$  mit  $R_{d'}(\theta) < \sigma^2/n$ . Wegen der Stetigkeit der Risikofunktion (in  $\theta$  für alle Entscheidungsregeln) gibt es also ein  $r > 0$ , so dass  $R_{d'}(\eta) < \sigma^2/n - r$  für  $|\eta - \theta| < r$ .

Sei nun  $\pi_b = \mathcal{N}(0, b^2)$  für  $b > 0$ . Einerseits ist damit  $r_{d'}(\pi_b) < \sigma^2/n$ . Andererseits ist nach Beispiel 3.18 der Bayes-Schätzer bezüglich  $\pi_b$  gerade

$$d_b(x) := \frac{nb^2}{nb^2 + \sigma^2} \bar{x}$$

mit Bayes-Risiko

$$\begin{aligned} r_{d_b}(\pi_b) &= \int \mathbb{E}_\theta[(d_b(X) - \theta)^2] \pi_b(d\theta) = \int \mathbb{V}_\theta[d_b(X)] + (\mathbb{E}_\theta[d_b(X)] - \theta)^2 \pi_b(d\theta) \\ &= \frac{n^2 b^4}{(nb^2 + \sigma^2)^2} \frac{\sigma^2}{n} + \frac{\sigma^4 b^2}{(nb^2 + \sigma^2)^2} = \frac{\sigma^2 n^2 b^4 + n \sigma^2 b^2}{n (nb^2 + \sigma^2)^2} \\ &= \frac{\sigma^2}{n} \left( \frac{n^2 b^4 + 2nb^2 \sigma^2 + \sigma^4 - nb^2 \sigma^2 - \sigma^4}{(nb^2 + \sigma^2)^2} \right) = \frac{\sigma^2}{n} \left( 1 - \frac{nb^2 \sigma^2 + \sigma^4}{(nb^2 + \sigma^2)^2} \right). \end{aligned}$$

Da  $d_b$  ein Bayes-Schätzer bezüglich  $\pi_b$  ist, gilt

$$\frac{2nb^2 \sigma^2 + \sigma^4}{(nb^2 + \sigma^2)^2} \geq \frac{\sigma^2}{n} - r_{d_b}(\pi_b) \geq \frac{\sigma^2}{n} - r_{d'}(\pi_b) \geq \frac{1}{\sqrt{2\pi} b^2} \int_{-\theta-r}^{\theta+r} r e^{-\eta^2/(2b^2)} d\eta.$$

Die linke Seite ist offenbar  $O(1/b^2)$  für  $b \rightarrow \infty$ , die rechte jedoch  $O(1/b)$ . Dies ist offenbar ein Widerspruch und damit ist die Zulässigkeit von  $d$  gezeigt. Nun ist  $d$  auch minimax nach Lemma 3.14.  $\square$

**Korollar 3.22 (Unbekannte Varianz).** *Betrachte dieselbe Situation wie in Proposition 3.21, jedoch mit unbekannter Varianz  $\sigma^2$ , d.h. das statistische Modell  $(X, \{\mathbb{P}_\theta : \theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+\})$  und  $g(\mu, \sigma^2) = \mu$ . Dann ist  $d(x) := \bar{x}$  zulässig und minimax-Schätzer.*

*Beweis.* Wie im Beweis von Proposition 3.21 ist nur die Zulässigkeit von  $d$  zu zeigen. Angenommen,  $d$  wäre unzulässig. Dann gibt es einen (oBdA nicht-randomisierten) Schätzer  $d'$  mit  $R_{d'} \leq R_d$  und ein  $\theta = (\mu, \sigma^2)$  mit  $R_{d'}(\theta) < R_d(\theta)$ . Da  $R_{d'}$  und  $R_d$  nicht davon abhängen, ob  $\sigma^2$  bekannt oder unbekannt ist, würde  $d'$  den Schätzer nun bei bekannter Varianz  $\sigma^2$  dominieren. Dies widerspricht aber der Aussage von Proposition 3.21 und die Behauptung ist gezeigt.  $\square$

**Bemerkung 3.23 (Höher-dimensionale Normalverteilungen).** Interessanterweise lässt sich Proposition 3.21 (und damit Korollar 3.22) nur bedingt auf höhere Dimensionen verallgemeinern. Ist nämlich in der gleichen Situation  $X = (X_1, \dots, X_n)$  mit  $X_i = (X_{i1}, \dots, X_{ik}) \in \mathbb{R}^k$  nach  $\mathbb{E}_\theta[X_{ij}] = \theta_j$  und  $\text{COV}[X_{ij}, X_{il}] = \delta_{jl}$ , so kann man zeigen, dass für die Verlustfunktion  $\ell(\theta, a) = \sum_{i=1}^k (\theta_i - a_i)^2 = (\theta - a)^\top (\theta - a)$  der Schätzer  $d(x) = \bar{x}$  für  $k \geq 3$  nicht mehr zulässig ist. Ein dominierender Schätzer ist durch den *James-Stein-Schätzer*

$$d'(x) := \left(1 - \frac{k-2}{x^\top x}\right) \bar{x}$$

gegeben. (Übrigens ist  $d'$  ebenfalls nicht zulässig und ist von

$$d''(x) := \left(1 - \frac{k-2}{x^\top x}\right)^+ \bar{x}$$

dominiert. Aber auch  $d''$  ist nicht zulässig.)

Nun zum Beweis der Unzulässigkeit von  $d$  im Fall  $k > 2$ . oBdA sei  $n = 1$ , also  $d(x) = x$ , andersfalls ändert sich nur die Varianz von  $d(X)$ . Wir berechnen die Risikofunktion von  $d$

$$R_d(\theta) \mathbb{E}_\theta[(\theta - d(X))^\top (\theta - d(X))] = k,$$

da  $(\theta - X)^\top (\theta - X) \sim \chi_k^2$ . Wir werden nun  $R_{d'}(\theta) < k$  zeigen. Es ist nämlich

$$\begin{aligned} R_{d'}(\theta) &= \mathbb{E}_\theta \left[ \left( X - \theta - \frac{(k-2)X}{X^\top X} \right)^\top \left( X - \theta - \frac{(k-2)X}{X^\top X} \right) \right] \\ &= \mathbb{E}_\theta[(X - \theta)^\top (X - \theta)] - 2(k-2) \mathbb{E}_\theta \left[ \frac{(X - \theta)^\top X}{X^\top X} \right] + (k-2)^2 \mathbb{E}_\theta \left[ \frac{1}{X^\top X} \right] \\ &= k - 2(k-2) \sum_{i=1}^k \mathbb{E}_\theta \left[ \frac{(X_i - \theta_i) X_i}{X^\top X} \right] + (k-2)^2 \mathbb{E}_\theta \left[ \frac{1}{X^\top X} \right]. \end{aligned}$$

Nun ist mit partieller Integration

$$\begin{aligned}
\sum_{i=1}^k \mathbb{E}_\theta \left[ \frac{(X_i - \theta_i) X_i}{X^\top X} \right] &= \sum_{i=1}^k \frac{1}{\sqrt{2\pi k}} \int \frac{x_1}{x^\top x} (x_i - \theta_i) \exp(- (x - \theta)^\top (x - \theta)/2) dx \\
&= \sum_{i=1}^k \frac{1}{\sqrt{2\pi k}} \int \frac{x^\top x - 2x_1^2}{(x^\top x)^2} \exp(- (x - \theta)^\top (x - \theta)/2) dx \\
&= \frac{k-2}{\sqrt{2\pi k}} \int \frac{1}{x^\top x} \exp(- (x - \theta)^\top (x - \theta)/2) dx \\
&= (k-2) \mathbb{E}_\theta \left[ \frac{1}{X^\top X} \right]
\end{aligned}$$

und damit

$$R_\theta(d') = k - (k-2) \mathbb{E}_\theta \left[ \frac{1}{X^\top X} \right] < k = R_\theta(d).$$

Wir wollen nun die entwickelte Theorie auf die Beispiele *Bern* und *Unif* anwenden.

**Beispiel 3.24 (Beispiel *Bern*).** Wir betrachten das Schätzproblem  $((X, \{\mathbb{P}_\theta : \theta \in (0, 1)\}, \aleph = (0, 1), g = \text{id}, \ell)$  für den Gauß-Verlust  $\ell$ . Speziell werden wir zeigen, dass  $d(x) = (\sum x_i)/n$  eine zulässige, aber keine minimax-Entscheidungsfunktion ist.

Wir wollen zunächst die Zulässigkeit von  $d$  zeigen. Offenbar minimiert  $d$  das Risiko genau dann für  $\ell(\theta, a) = (\theta - a)^2$ , wenn  $d$  das Risiko für  $\ell'(\theta, a) = (\theta - a)^2/(\theta(1 - \theta))$  minimiert. Sei nun  $\pi = U(0, 1)$  eine a-priori-Verteilung. Da dies genau die  $\beta(1, 1)$ -Verteilung ist, sehen wir aus Beispiel 2.38, dass  $\pi_x = \beta(1 + \sum x_i, n + 1 - \sum x_i)$  die a-posteriori-Verteilung ist. Um den Bayes-Schätzer zu bestimmen, ist hierfür nach Theorem 3.16.2

$$\mathbb{E}[\ell(\Theta_x, a)] = \frac{\Gamma(n+2)}{\Gamma(1 + \sum x_i) \Gamma(n+1 - \sum x_i)} \int_0^1 (\theta - a)^2 \theta^{\sum x_i - 1} (1 - \theta)^{n-1 - \sum x_i} d\theta$$

zu minimieren. Offenbar liegt dieses Minimum genau beim Erwartungswert der  $\beta(\sum x_i, n - \sum x_i)$ -Verteilung, also bei<sup>8</sup>  $a = (\sum x_i)/n$ . Damit ist  $d$  Bayes-Schätzer und nach Lemma 3.19 auch zulässig.

Um zu zeigen, dass  $d$  nicht minimax ist sei  $d_{a,b}(x) := \frac{a + \sum x_i}{a + b + n}$ . Für die a-priori-Verteilung  $\pi = \beta(a, b)$  ist die a-posteriori-Verteilung  $\beta(a + \sum x_i, b + n - \sum x_i)$  und der Bayes-Schätzer ist gerade  $d_{a,b}(x)$ . Das Risiko dieses Schätzers ist

$$\begin{aligned}
R_{d_{a,b}}(\theta) &= \mathbb{E}_\theta[(d_{a,b}(X) - \theta)^2] = \mathbb{V}_\theta[d_{a,b}(X)] + (\mathbb{E}_\theta[d_{a,b}(X)] - \theta)^2 \\
&= \frac{1}{(a + b + n)^2} (n\theta(1 - \theta) + (a + n\theta - (a + b + n)\theta)^2) \\
&= \frac{1}{(a + b + n)^2} (\theta^2((a + b)^2 - n) + \theta(n - 2a(a + b)) + a^2).
\end{aligned}$$

Für  $a = b = \sqrt{n}/2$  ist dieser konstant  $\frac{n/4}{(n + \sqrt{n})^2}$ , und damit ist  $d_{\sqrt{n}/2, \sqrt{n}/2}$  nach Theorem 3.20 minimax. Allerdings ist

$$\sup_{\theta \in (0,1)} R_d(\theta) = \sup_{\theta \in (0,1)} R_{d_{0,0}}(\theta) = \sup_{\theta \in (0,1)} \frac{n\theta(1 - \theta)}{n^2} = \frac{1}{4n} > \frac{n/4}{(n + \sqrt{n})^2}$$

und damit ist  $d$  nicht minimax.

<sup>8</sup>Wir verwenden, dass für  $\beta(p, q)$  der Erwartungswert durch  $p/(p + q)$  und die Varianz durch  $pq/((p + q + 1)(p + q)^2)$  gegeben ist.

## 4 Testtheorie

In diesem Kapitel sei immer  $((X, \{\mathbb{P}_\theta : \theta \in \Theta\}), \mathfrak{N}, \ell)$  ein statistisches Entscheidungsproblem mit  $\mathfrak{N} = \{H_0, H_1\}$  (d.h. die Entscheidung besteht immer zwischen der Null-Hypothese  $H_0$  und der Alternativ-Hypothese  $H_1$ ). Wir erinnern an die Neyman-Pearson'schen Verlustfunktion

$$\ell(\theta, H_1) = \begin{cases} 0, & \theta \in \mathcal{P}_1, \\ \ell_0, & \theta \in \mathcal{P}_0 \end{cases},$$

$$\ell(\theta, H_0) = \begin{cases} 0, & \theta \in \mathcal{P}_0, \\ \ell_1, & \theta \in \mathcal{P}_1 \end{cases}.$$

für  $\mathcal{P}_0, \mathcal{P}_1$  mit  $\mathcal{P} = \mathcal{P}_0 \uplus \mathcal{P}_1$ . Wir werden immer  $\ell_0 = \ell_1 = 1$  betrachten. (Der allgemeinere Fall ergibt sich dann meist ebenfalls.) Die Entscheidung  $a = H_0$  bedeutet, dass man sich für die Nullhypothese entscheidet (bzw. dass man sie nicht verwirft) und  $a = H_1$  bedeutet, dass man die Nullhypothese verwirft und sich für die Alternative entscheidet. Eine Entscheidungsfunktion (bzw. der Test)  $\delta$  wird eindeutig durch die Funktion

$$\varphi(x) := \delta(x, \{H_1\})$$

definiert, d.h. durch die Wahrscheinlichkeit, sich für die Alternativ-Hypothese zu entscheiden bei Vorliegen der Daten  $x$ . (Damit ist automatisch  $\delta(x, \{H_0\}) = 1 - \varphi(x)$ .) Das Entscheidungsproblem (oder auch Test-Problem) heißt *einfach*, wenn sowohl  $\mathcal{P}_0$  als auch  $\mathcal{P}_1$  ein-elementig sind. Andernfalls heißt es *zusammengesetzt*. Wir erinnern an die Gütefunktion  $\beta_\varphi : \theta \mapsto \mathbb{E}_\theta[\varphi(X)]$ . Das Risiko von  $\delta$  (oder  $\varphi$ ) ist gegeben durch

$$R_\varphi(\theta) := R_\delta(\theta) = \mathbb{E}_\theta \left[ \int \ell(\theta, a) \delta(X, da) \right] = \begin{cases} \mathbb{E}_\theta[\varphi(X)] = \beta_\varphi(\theta), & \theta \in \mathcal{P}_0, \\ 1 - \mathbb{E}_\theta[\varphi(X)] = 1 - \beta_\varphi(\theta), & \theta \in \mathcal{P}_1. \end{cases}$$

Weiter sei  $\Phi$  die Menge aller möglichen Tests (gegeben entweder durch  $\delta$  oder  $\varphi$ ) und  $\{\theta \mapsto R_\delta(\theta) : \delta \in \Phi\}$  die Risikomenge, d.h. die Menge der Risikofunktionen aller Tests. Für einfache Tests mit  $\mathcal{P}_0 = \{\mathbb{P}_0\}, \mathcal{P}_1 = \{\mathbb{P}_1\}$  identifizieren wir

$$\mathcal{R} = \{(\mathbb{E}_0[\varphi(X)], 1 - \mathbb{E}_1[\varphi(X)]) : \varphi \in \Phi\}. \quad (\text{R})$$

### 4.1 Bayes-Tests

Wie wir in Theorem 3.16 gesehen haben, ist es oftmals gar nicht schwierig, Bayes-Entscheidungsfunktionen aufzustellen. Dies wollen wir nun für Tests durchführen, die entsprechenden Entscheidungsfunktionen heißen dann Bayes-Tests.

**Proposition 4.1 (Bayes-Tests).** *Sei  $\pi$  ein Wahrscheinlichkeitsmaß auf  $\mathcal{P}$ . Dann ist  $\varphi \in \Phi$  genau dann ein Bayes-Test (d.h. eine Bayes-Entscheidungsfunktion für das Test-Problem) bezüglich  $\pi$ , wenn (für die a-posteriori-Verteilung  $\pi_x$ )*

$$\varphi(x) = \begin{cases} 1, & \pi_x(\mathcal{P}_0) < \pi_x(\mathcal{P}_1), \\ 0, & \pi_x(\mathcal{P}_0) > \pi_x(\mathcal{P}_1). \end{cases}$$

(Auf dem Bereich  $\{x : \pi_x(\mathcal{P}_0) = \pi_x(\mathcal{P}_1)\}$  ist  $\varphi$  nicht eindeutig bestimmt.)

*Beweis.* Wir berechnen das Bayes-Risiko für  $\Theta \sim \pi$  für  $\Theta_x \sim \pi_x$  und den Test  $\psi \in \Phi$  als

$$\begin{aligned} r_\delta(\pi) &= \mathbb{E} \left[ \mathbb{E} \left[ \int \ell(\Theta, a) \delta(X, da) \middle| X \right] \right] = \mathbb{E} \left[ \mathbb{E} \left[ \ell(\Theta, H_1) \psi(X) + \ell(\Theta, H_0) (1 - \psi(X)) \middle| X \right] \right] \\ &= \mathbb{E} [\mathbb{E} [\mathbb{1}_{\Theta_X \in \mathcal{P}_0} \psi(X) + \mathbb{1}_{\Theta_X \in \mathcal{P}_1} (1 - \psi(X)) | X]] \\ &= \mathbb{E} [\pi_X(\mathcal{P}_1) + (\pi_X(\mathcal{P}_0) - \pi_X(\mathcal{P}_1)) \psi(X)]. \end{aligned}$$

Die Funktion  $\psi$ , die die rechte Seite minimiert, muss genau die angegebene Form haben. Daraus folgt die Behauptung.  $\square$

**Korollar 4.2 (Bayes-Tests für einfache Tests).** Sei  $\mathcal{P}_0 = \{\mathbb{P}_0 = p_0 \cdot \lambda\}$ ,  $\mathcal{P}_1 = \{\mathbb{P}_1 = p_1 \cdot \lambda^n\}$  (d.h. wir betrachten einen einfachen Test für ein reguläres statistisches Modell) und  $\pi_k(\mathcal{P}_0) = \frac{k}{k+1}$  (und damit  $(\pi(\mathcal{P}_1) = \frac{1}{k+1})$ ) für ein  $k \in [0, \infty]$ <sup>9</sup>. Dann ist  $\varphi_k$  genau dann ein Bayes-Test bezüglich  $\pi_k$ , wenn  $\varphi_k$  gegeben ist als

$$\begin{aligned} \varphi_k(x) &:= \begin{cases} 1, & \frac{p_1(x)}{p_0(x)} > k, \\ 0, & \frac{p_1(x)}{p_0(x)} < k \end{cases}, \\ \varphi_0(x) &:= \mathbb{1}_{p_1(x) > 0}, \\ \varphi_\infty(x) &:= \mathbb{1}_{p_0(x) > 0}. \end{aligned}$$

*Beweis.* Wir zeigen nur den Fall  $0 < k < \infty$ , die anderen beiden Fälle zeigt man direkt. Es gilt

$$\pi_x(\{H_0\}) = \frac{p_0(x)k}{p_0(x)k + p_1(x)}, \quad \pi_x(\{H_1\}) = \frac{p_1(x)}{p_0(x)k + p_1(x)}.$$

Aus Proposition 4.1 folgt, dass  $\varphi_k$  Bayes-Tests sind mit

$$\pi_x(\{\mathbb{P}_0\}) < \pi_x(\{\mathbb{P}_1\}) \text{ genau dann, wenn } \frac{p_1(x)}{p_0(x)} > k.$$

$\square$

## 4.2 Likelihood-Quotienten-Tests

Die soeben konstruierten Bayes-Tests werden gerade durch die Quotienten der Dichten der Alternative und der Nullhypothese bestimmt. Diese Dichten sind – bei gegebenen Daten – außerdem als Likelihood (des Parameters) bekannt. Deshalb nennt man diese Tests auch Likelihood-Quotienten-Tests.

**Definition 4.3 (Likelihood-Quotienten-Test).** Sei  $\mathcal{P}_0 = \{\mathbb{P}_0 = p_0 \cdot \lambda^n\}$ ,  $\mathcal{P}_1 = \{\mathbb{P}_1 = p_1 \cdot \lambda^n\}$ . Für  $k \in [0, \infty]$  heißt  $\varphi_k$  aus Korollar 4.2 Likelihood-Quotienten-Test (oder LQ-Test) mit kritischem Wert  $k$ . Genauer setzen wir

$$\varphi_k(x) := \varphi_{k,\gamma}(x) := \begin{cases} 1, & \frac{p_1(x)}{p_0(x)} > k, \\ \gamma(x), & \frac{p_1(x)}{p_0(x)} = k, \\ 0, & \frac{p_1(x)}{p_0(x)} < k. \end{cases}$$

Die Menge  $\{x : \frac{p_1(x)}{p_0(x)} = k\}$  heißt Randomisierungsbereich von  $\varphi_k$ .

<sup>9</sup>Wir setzen  $\infty/\infty = 1$ .

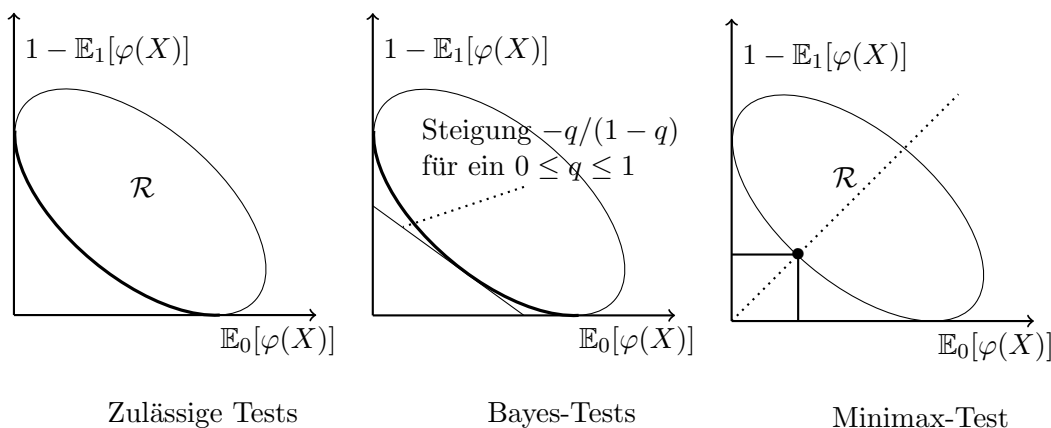
**Beispiel 4.4 (Beispiel *Bern*).** Im Beispiel *Bern* seien  $\theta_0, \theta_1 \in (0, 1)$  verschieden. Ein Likelihood-Quotiententest ist von der Form

$$\varphi(x) = \begin{cases} 1, & \frac{\theta_1^{\sum x_i} (1-\theta_1)^{n-\sum x_i}}{\theta_0^{\sum x_i} (1-\theta_0)^{n-\sum x_i}} > k, \\ 0, & \frac{\theta_1^{\sum x_i} (1-\theta_1)^{n-\sum x_i}}{\theta_0^{\sum x_i} (1-\theta_0)^{n-\sum x_i}} < k, \end{cases}$$

also für ein geeignetes  $k'$

$$\varphi(x) = \begin{cases} 1, & \left(\frac{\theta_1(1-\theta_0)}{\theta_0(1-\theta_1)}\right)^{\sum x_i} > k', \\ 0, & \left(\frac{\theta_1(1-\theta_0)}{\theta_0(1-\theta_1)}\right)^{\sum x_i} < k', \end{cases}$$

**Bemerkung 4.5 (Grafische Darstellung der Risikomenge für einfache Tests).** Die Menge aller Tests  $\Phi$  ist konvex, da die Konvexkombination zweier  $[0, 1]$ -wertiger Funktionen wieder eine solche Funktion ergibt. Demnach ist auch die in (R) definierte Risikomenge für einfache Tests eine konvexe Teilmenge des  $\mathbb{R}_+^2$ . Weiter ist sowohl  $\varphi = 0$  als auch  $\varphi = 1$  erlaubt, so dass  $\mathcal{R}$  sowohl die  $x$ - als auch die  $y$ -Achse berührt. Hieraus lassen sich die zulässigen, Bayes-Tests (also LQ-Tests) und minimax-Tests ablesen.



*Zulässige Tests* ergeben sich für solche  $\varphi$ , für die zwar  $R_\varphi \in \mathcal{R}$ , aber kein  $(x, y) \leq R_\varphi$  (in der üblichen Halbordnung in  $\mathbb{R}_+^2$ ) in  $\mathcal{R}$  ist.

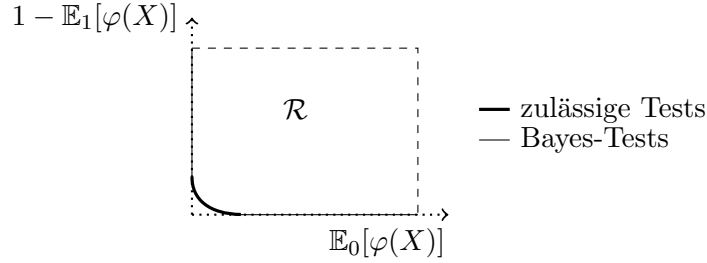
Ein *Bayes-Test*  $\varphi$  bezüglich  $\pi = (q, 1 - q)$  minimiert gerade das Skalarprodukt

$$r_\varphi(\pi) = (q, 1 - q)(\mathbb{E}_0[\varphi(X)], 1 - \mathbb{E}_1[\varphi(X)]) = \min_{(x, y) \in \mathcal{R}} (q, 1 - q)(x, y).$$

Schreibt man  $(x, y) = a(0, 1) + b(1 - q, -q)$  (d.h. eine parallelverschobene Gerade mit Steigung  $-q/(1 - q)$ ), so ist  $(q, 1 - q)(x, y) = a(1 - q)$ . Dieses Minimum ergibt sich also gerade als die am wenigsten verschobene Gerade mit Steigung  $-q/(1 - q)$ , die  $\mathcal{R}$  berührt.

Wir wissen aus Theorem 3.20, dass ein zulässiger Test mit konstantem Risiko ein *minimax-Test* ist. Dies ist gerade für  $\mathbb{E}_0[\varphi(X)] = 1 - \mathbb{E}_1[\varphi(X)]$  gegeben, also für Schnittpunkte von  $\mathcal{R}$  mit  $\{(x, x) : x \geq 0\}$ .

Es gibt in diesem Bild die Möglichkeit, dass die Menge der Bayes-Tests größer ist als die der zulässigen Tests, was wir nun veranschaulichen wollen.



Hier bestehen die minimalen Geraden, die  $\mathcal{R}$  berühren, ebenfalls aus den horizontalen und vertikalen Begrenzungen von  $\mathcal{R}$ . Diese sind jedoch keine zulässigen Tests, weil es Punkte  $(x, y) \in \mathcal{R}$  gibt, die diese Tests dominieren.

**Proposition 4.6 (Zulässige und minimax-LQ-Tests).** Sei  $\mathcal{P}_0 = \{\mathbb{P}_0 = p_0 \cdot \lambda^n\}$ ,  $\mathcal{P}_1 = \{\mathbb{P}_1 = p_1 \cdot \lambda^n\}$  und  $\varphi_k$  für  $k \in [0, \infty]$  ein Likelihood-Quotienten-Test mit kritischem Parameter  $k$ . Dann gilt:

1. Ist  $\varphi$  ein zulässiger Test, so gibt es  $k \in [0, \infty]$  mit  $\varphi = \varphi_k$ .
2. Gilt  $\mathbb{E}_0[\varphi_k(X)] > 0$ ,  $\mathbb{E}_1[\varphi_k(X)] < 1$ , so ist  $\varphi_k$  zulässig.
3. Genau dann ist  $\varphi$  ein minimax-Test, wenn es  $k \in [0, \infty]$  gibt mit  $\varphi = \varphi_k$  und  $\mathbb{E}_0[\varphi_k(X)] = \mathbb{E}_1[1 - \varphi_k(X)]$ .

*Beweis.* 1. Wir verwenden die grafische Darstellung aus Bemerkung 4.5. Ist  $\varphi$  mit  $(\mathbb{E}_0[\varphi(X), 1 - \mathbb{E}_1[\varphi(X)]) \in \mathcal{R}$  zulässig (also Element des unteren Randes der Risikomenge), so lesen wir ab, dass es eine minimale berührende Gerade gibt, die  $\mathcal{R}$  gerade in  $(\mathbb{E}_0[\varphi(X), 1 - \mathbb{E}_1[\varphi(X)])$  berührt. Damit ist  $\varphi$  ein Bayes-Test, also nach Korollar 4.2 auch ein LQ-Test.

2. Wir müssen ausschließen, dass  $k = 0$  oder  $k = \infty$ . Für  $k \in (0, \infty)$  ist nämlich  $\varphi_k$  Bayes-Test zu einer a-priori-Verteilung mit vollem Träger und damit nach Lemma 3.19 zulässig. Offenbar ist  $\mathbb{P}_i[p_i(X) > 0] = 1$ ,  $i = 0, 1$ . Wäre  $k = 0$ , so gilt aber  $\mathbb{E}_1[\varphi_0(X)] = \mathbb{P}_1(p_1(X) > 0) = 1$  im Widerspruch zur Voraussetzung und wäre  $k = \infty$ , so gilt  $\mathbb{E}_0[\varphi_\infty(X)] = \mathbb{P}_0[p_0(X) > 0] = 1$ .

3. ' $\Leftarrow$ ': Offenbar ist  $\varphi_k$  Bayes-Test mit konstantem Risiko. Deshalb folgt die Behauptung aus Theorem 3.20.2.

' $\Rightarrow$ ': Wir zeigen zunächst, dass  $\mathbb{E}_0[\varphi(X)] = 1 - \mathbb{E}_1[\varphi(X)]$  gelten muss. Angenommen, es wäre  $\lambda := 1 - \mathbb{E}_1[\varphi(X)] - \mathbb{E}_0[\varphi(X)] > 0$ . Wir definieren dann die neue Entscheidungsfunktion

$$\psi(X) := \frac{1}{1 + \lambda} \varphi(X) + \frac{\lambda}{1 + \lambda} \in \Phi.$$

Nun gilt

$$\begin{aligned} \mathbb{E}_0[\psi(X)] &= \frac{1}{1 + \lambda} \mathbb{E}_0[\varphi(X)] + \frac{1}{1 + \lambda} (1 - \mathbb{E}_1[\varphi(X)] - \mathbb{E}_0[\varphi(X)]) \\ &= \frac{1}{1 + \lambda} (1 - \mathbb{E}_1[\varphi(X)]) \\ &= 1 - \frac{1}{1 + \lambda} \mathbb{E}_1[\varphi(X)] - \frac{\lambda}{1 + \lambda} = 1 - \mathbb{E}_1[\psi(X)]. \end{aligned}$$

Weiter gilt

$$\begin{aligned} \max(\mathbb{E}_0[\psi(X)], 1 - \mathbb{E}_1[\psi(X)]) &= \frac{1}{1 + \lambda}(1 - \mathbb{E}_1[\varphi(X)]) < 1 - \mathbb{E}_1[\varphi(X)] \\ &= \max(\mathbb{E}_0[\varphi(X)], 1 - \mathbb{E}_1[\varphi(X)]). \end{aligned}$$

Dies steht im Widerspruch dazu, dass  $\varphi$  minimax ist. Im Fall  $\lambda < 0$  führt eine ähnliche Konstruktion zum Ziel. Nun ist also  $(\mathbb{E}_0[\varphi(X)], 1 - \mathbb{E}_0[\varphi(X)]) \in \mathcal{R} \cap \{(x, x) : x \in \mathbb{R}_+\}$ . Aus der grafischen Darstellung aus Bemerkung 4.5 und der Konvexität der Risikomenge  $\mathcal{R}$  folgt, dass  $\varphi$  ein Bayes-Test sein muss und damit ein LQ-Test.  $\square$

Auch für zusammengesetzte Hypothesen kann man Likelihood-Quotiententests definieren. Dies wollen wir nun tun, auch wenn wir in diesem Fall nur Beispiele angeben und die Optimalität der Tests nicht weiter untersuchen werden.

**Definition 4.7 (Likelihood-Quotienten-Tests für zusammengesetzte Alternativen).** Sei  $(X, \{\mathcal{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}_0\}, \mathfrak{N}, \ell)$  ein reguläres statistisches Modell,  $\mathcal{P} = \mathcal{P}_0 \uplus \mathcal{P}_1$  eine Partition von  $\mathcal{P}$  und  $H_0 : \theta \in \mathcal{P}_0$ ,  $H_1 : \theta \in \mathcal{P}_1$ , sowie  $\mathfrak{N} = \{H_0, H_1\}$  und  $\ell$  der Neyman-Pearson-Verlust. Dann heißt ein Test  $\varphi_k$  mit

$$\varphi_k(x) := \varphi_{k,\gamma}(x) := \begin{cases} 1, & \lambda(x) := \frac{\sup_{\theta \in \mathcal{P}_0} p_\theta(x)}{\sup_{\theta \in \mathcal{P}} p_\theta(x)} < k, \\ \gamma(x), & \lambda(x) = k, \\ 0, & \lambda(x) > k \end{cases}$$

Likelihood-Quotienten-Test mit kritischem Wert  $k$ .

**Bemerkung 4.8 (Mögliche Werte für  $k$ ).** Anders als bei Likelihood-Quotienten-Tests von einfachen Hypothesen ist die Definition bei zusammengesetzten Hypothesen so, dass immer  $\lambda \leq 1$  ist. Es machen also nur  $k \in [0, 1]$  Sinn. Diese Werte legen natürlich auch das Signifikanzniveau fest.

**Beispiel 4.9 (Test auf den Parameter einer Exponentialverteilung).** Sei  $\mathcal{P} = [\theta_0, \infty)$  und  $\mathbb{P}_\theta = \exp(\theta)^n$ . Es ist unter  $\mathbb{P}_\theta$  also  $X = (X_1, \dots, X_n)$  unabhängig und  $X_i \sim \exp(\theta)$ . Weiter sei  $\mathcal{P}_0 = \{\theta_0\}$  und  $\mathcal{P}_1 = (\theta_0, \infty)$ . Nun ist  $\varphi$  genau dann ein Likelihood-Quotienten-Test, falls

$$\varphi(x) := 1_{\{\sum_{i=1}^n x_i > c\}}$$

für ein  $c > 0$ .

Denn: Zunächst berechnen wir  $\sup_{\theta \in \mathcal{P}} p_\theta(x)$ . Wir schreiben hierfür

$$\log p_\theta(x) = \log(\theta^n e^{-\theta(x_1 + \dots + x_n)}) = n \log \theta - \theta(x_1 + \dots + x_n)$$

und damit  $\frac{d}{d\theta} \log p_\theta(x) = \frac{n}{\theta} - (x_1 + \dots + x_n)$ , also

$$\log \frac{p_{\theta_0}(x)}{\sup_{\theta \in \mathcal{P}} p_\theta(x)} = \begin{cases} \log p_{\theta_0}(x) - \log p_{1/\bar{x}}(x) = n \log \theta_0 - \theta_0 n \bar{x} + n \log \bar{x} - n, & \bar{x} < 1/\theta_0, \\ 0, & \bar{x} \geq 1/\theta_0. \end{cases}$$

Damit ist ein Likelihood-Quotienten-Test von der Form  $\varphi(x) = 1_{\log \bar{x} - \theta_0 \bar{x} < k} 1_{\bar{x} > 1/\theta_0}$ . Nun ist  $\bar{x} \mapsto \log \bar{x} - \theta_0 \bar{x}$  für  $\bar{x} \geq 1/\theta_0$  monoton fallend, d.h.  $\varphi$  ist von der Form  $\varphi(x) = 1_{\bar{x} > k'}$ .



Bereits bekannte Tests, etwa der  $t$ -Test, sind ebenfalls Likelihood-Quotienten-Tests.

**Proposition 4.10 (Einfacher  $t$ -Test ist ein Likelihood-Quotienten-Test).** Sei  $(X, \{\mathbb{P}_\theta = \mathcal{N}(\mu, \sigma^2)^n : \theta = (\mu, \sigma)^2 \in \mathcal{P}_0 := \mathbb{R} \times \mathbb{R}_+\})$  das Normalverteilungsmodell,  $\mathcal{P}_0 = \{\mathbb{P}_\theta : \theta = (\mu_0, \sigma^2) \in \{\mu_0\} \times \mathbb{R}_+\}$ ,  $\mathcal{P}_1 = \mathcal{P} \setminus \mathcal{P}_0$ ,  $\mathfrak{K} = \{H_0, H_1\}$  und  $\ell$  der Neyman-Pearson-Verlust. Dann ist  $\varphi$  genau dann ein Likelihood-Quotienten-Test, wenn  $\varphi(x) = 1_{|t(x)| \geq c}$  für ein geeignetes  $c$  mit

$$t(x) = \frac{\bar{x} - \mu_0}{\sqrt{\hat{s}^2(x)/n}}.$$

*Beweis.* Sei  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  mit

$$p_\theta(x) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\sum_{i=1}^n (x_i - \theta)^2}{2\sigma^2}\right).$$

Bekannt ist, dass der Maximum-Likelihood-Schätzer für  $\theta = (\mu, \sigma^2)$  immer durch  $(\bar{x}, (n-1)\hat{s}^2(x)/n)$  gegeben ist. Bei gegebenen  $\mu$  ist der Maximum-Likelihood-Schätzer von  $\sigma^2$  durch  $\hat{s}^2(x) := \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$  gegeben. Um also  $p_\theta$  zu maximieren, berechnen wir erst

$$\begin{aligned} \sup_{\sigma^2 \in \mathbb{R}_+} p_\theta(x) &= (2\pi\hat{s}^2(x))^{-n/2} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\hat{s}^2(x)}\right) \\ &= \left(\frac{2\pi}{n} \sum_{i=1}^n (x_i - \mu)^2\right)^{-n/2} e^{-n/2} \end{aligned}$$

und damit

$$\begin{aligned} \lambda(x) &= \left(\frac{\sum_{i=1}^n (x_i - \mu_0)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)^{-n/2} = \left(\frac{\sum_{i=1}^n (x_i - \bar{x})^2 + (\bar{x} - \mu_0)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)^{-n/2} \\ &= \left(1 + |t(x)|^2 \frac{n-1}{n}\right)^{-n/2}. \end{aligned}$$

Die Behauptung folgt nun daraus, dass  $x \mapsto (1 + (n-1)x/n)^{-n/2}$  monoton ist.  $\square$

**Bemerkung 4.11 (Approximative Verteilung des Likelihood-Quotienten).** Um das Signifikanzniveau eines Tests  $\varphi$  zu bestimmen, benötigt man  $\sup_{\theta \in \mathcal{P}_0} \mathbb{E}_\theta[\varphi(X)]$ . Da bei Likelihood-Quotienten-Tests  $\varphi = \varphi_{k,\gamma}$  von den Parametern  $k$  und  $\gamma$  abhängt, heißt das, dass man diese so bestimmen muss, dass

$$\mathbb{E}_\theta[\varphi_{k,\gamma}(X)] = \mathbb{P}_\theta[\lambda(X) \leq k] + \mathbb{E}_\theta[\gamma(X), \lambda(X) = k] \leq \alpha, \quad \theta \in \mathcal{P}_0.$$

Hierzu benötigt man also die Verteilung des Likelihood-Quotienten  $\lambda$ .

Sei  $\mathcal{P} \subseteq \mathbb{R}^p$  und  $\mathcal{P}_0 = \{\eta\}$ . Wir werden (mit Hilfe der noch zu zeigenden Aussagen aus Abschnitt 5.4) nun zeigen, dass (unter gewissen Regularitätsannahmen)

$$-2 \log \lambda(X) =: \Lambda(X) \xrightarrow{n \rightarrow \infty} Z \sim \chi_p^2.$$

Sei  $\hat{\theta}$  der Maximum-Likelihood-Schätzer für  $\theta$ . Dann schreiben wir

$$\begin{aligned} \sum_{k=1}^n \log p_\eta^n(x_k) &= \sum_{k=1}^n \left( \log p_{\hat{\theta}}(x_k) + \sum_{i=1}^p \frac{\partial}{\partial \theta_i} \log p_{\hat{\theta}}(x_k) (\eta_i - \hat{\theta}_i) \right. \\ &\quad \left. + \frac{1}{2} \sum_{i,j=1}^k \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \log p_{\hat{\theta}}(x_k) (\eta_i - \hat{\theta}_i) (\eta_j - \hat{\theta}_j) \right) + o((\eta - \hat{\theta})^2). \end{aligned}$$

Nun verschwindet der zweite Term, da  $\hat{\theta}$  die Maximalstelle von  $\log p_{\theta}$  ist, und damit ist für  $X \sim \mathbb{P}_{\eta}^n$

$$\begin{aligned}\Lambda(X) &= 2 \sum_{k=1}^n \log p_{\hat{\theta}}(X_k) - 2 \log p_{\eta}(X_k) \\ &= - \sum_{k=1}^n \sum_{i,j=1}^p \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \log p_{\hat{\theta}}(X_k) (\eta_i - \hat{\theta}_i) (\eta_j - \hat{\theta}_j) + o((\eta - \hat{\theta})^2) \\ &\stackrel{n \rightarrow \infty}{\approx} -n \sum_{i,j=1}^p \mathbb{E}_{\eta} \left[ \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \log p_{\eta}(X) \right] (\eta_i - \hat{\theta}_i) (\eta_j - \hat{\theta}_j),\end{aligned}$$

da  $\hat{\theta}$  konsistent ist (siehe Theorem 5.30). Nun ist nach Theorem 5.31 und Bemerkung 5.18 für

$$I(\theta) := - \left( \mathbb{E}_{\theta} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p_{\theta}(X) \right] \right)_{i,j=1,\dots,k}$$

gerade (für die Einheitsmatrix  $E_k$ )

$$\sqrt{n} (I(\eta))^{1/2} (\hat{\theta} - \eta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, E_p),$$

also

$$\Lambda(X) \approx n (\hat{\theta} - \eta) I(\eta) (\hat{\theta} - \eta) \xrightarrow{n \rightarrow \infty} Z_1^2 + \dots + Z_p^2 \sim \chi_p^2$$

### 4.3 Beste Tests

Optimale Tests minimieren die Risikofunktion. Allerdings ist es nicht sinnvoll, alle Tests zuzulassen, sondern nur solche, deren Fehler erster Art kleiner als ein vorgegebenes  $\alpha$  sind. Diese Tests heißen auch beste Tests.

**Definition 4.12 (Niveau eines Tests).** 1. Für  $\alpha \in [0, 1]$  ist

$$\Phi_{\alpha} := \{\varphi \in \Phi : \mathbb{E}_{\theta}[\varphi(X)] \leq \alpha \text{ für alle } \theta \in \mathcal{P}_0\}$$

die Menge aller Tests zum Niveau  $\alpha$ .

2. Ein Test  $\varphi$  heißt (gleichmäßig) bester Tests zum Niveau  $\alpha$  (oder (Uniformly) Most Powerful Test oder UMP-Test), falls

$$\mathbb{E}_{\theta}[\varphi(X)] = \sup_{\psi \in \Phi_{\alpha}} \mathbb{E}_{\theta}[\psi(X)], \quad \theta \in \mathcal{P}_1.$$

Das folgende Resultat erlaubt die Konstruktion bester Tests im Falle von einfachen Hypothesen.

**Theorem 4.13 (Neyman-Pearson-Lemma).** Wir betrachten das reguläre statistische Modell  $(X, \{\mathbb{P}_{\theta} = p_{\theta} \cdot \lambda^n : \theta \in \mathcal{P}\})$  mit  $\mathcal{P}_i = \{\mathbb{P}_i\}$ ,  $i = 0, 1$  und  $\mathcal{P} = \mathcal{P}_0 \cup \mathcal{P}_1$ . Sei  $\alpha \in (0, 1)$ .

1. Es gibt einen LQ-Test  $\varphi_{k,\gamma}$  mit  $\gamma \in [0, 1]$  konstant und  $\mathbb{E}_0[\varphi(X)] = \alpha$ .
2. Ist  $\varphi$  ein LQ-Test mit  $\mathbb{E}_0[\varphi(X)] = \alpha$ , so ist  $\varphi$  bester Test zum Niveau  $\alpha$ .

3. Ist  $\varphi$  ein bester Test zum Niveau  $\alpha$ , so ist  $\varphi$  ein LQ-Test. Es gilt dann entweder  $\mathbb{E}_0[\varphi(X)] = \alpha$  oder  $\mathbb{E}_1[\varphi(X)] = 1$ .

*Beweis.* 1. Wir setzen

$$k := \inf\{y : \mathbb{P}_0(p_1(X)/p_0(X) > y) \leq \alpha\}$$

und haben damit  $k$  bereits bestimmt. Es gilt  $\mathbb{P}_0(p_1(X)/p_0(X) \geq k) \geq \alpha \geq \mathbb{P}_0(p_1(X)/p_0(X) > k)$ . Für die Wahl von  $\gamma$  gibt es nun zwei Möglichkeiten. Ist  $\mathbb{P}_0(p_1(X)/p_0(X) = k) = 0$ , so setzen wir  $\gamma = 0$ . Andernfalls setzen wir

$$\gamma := \frac{\alpha - \mathbb{P}_0(p_1(X)/p_0(X) > k)}{\mathbb{P}_0(p_1(X)/p_0(X) = k)} \in (0, \infty).$$

Nun gilt nach Definition des LQ-Tests  $\varphi_{k,\gamma}^*$

$$\mathbb{E}_0[\varphi_{k,\gamma}^*(X)] = \mathbb{P}_0(p_1(X)/p_0(X) > k) + \gamma \mathbb{P}_0(p_1(X)/p_0(X) = k) = \alpha.$$

2. Sei  $\varphi = \varphi_k^*$  für ein  $k \in [0, \infty]$  sowie  $\psi \in \Phi_\alpha$ . Wir müssen zeigen, dass

$$\mathbb{E}_1[\varphi_k^*(X)] \geq \mathbb{E}_1[\psi(X)]$$

und schreiben hierfür

$$\begin{aligned} & \mathbb{E}_1[\varphi^*(X)] - \mathbb{E}_1[\psi(X)] \\ &= \int (\varphi^*(x) - \psi(x))(p_1(x) - kp_0(x))\lambda^n(dx) + k \int (\varphi_k^*(x) - \psi(x))p_0(x)\lambda^n(dx) \\ &= \int (1_{p_1(x) > kp_0(x)} - \psi(x))(p_1(x) - kp_0(x))\lambda^n(dx) + k\mathbb{E}_0[\varphi_k^*(X) - \psi(X)] \\ &\geq \int 1_{p_1(x) > kp_0(x)}(1 - \psi(x))(p_1(x) - kp_0(x))\lambda^n(dx) \geq 0. \end{aligned}$$

3. Sei  $\varphi$  bester Test zum Niveau  $\alpha$  und  $\varphi_k^*$  der LQ-Test zum Niveau  $\alpha$  aus 1. Nach 2. ist  $\varphi_k^*$  ebenfalls bester Test zum Niveau  $\alpha$  und wir müssen zeigen, dass  $\varphi = \varphi_k^*$ . In der Rechnung aus dem Beweis von 2. muss Gleichheit in beiden Abschätzungen gelten, es gilt also  $\mathbb{E}_0[\varphi(X)] = \mathbb{E}_0[\varphi_k^*(X)] = \alpha$  sowie  $(\varphi_k^* - \varphi)(p_1 - kp_0) \stackrel{\lambda^n\text{-f.ü.}}{=} 0$ . Daraus folgt  $\{x : \varphi_k^*(x) \neq \varphi(x)\} \subseteq \{x : p_1(x)/p_0(x) = k\}$ , d.h.  $\varphi$  muss ein LQ-Test sein, der höchstens auf dem Randomisierungsbereich nicht mit  $\varphi_k^*$  übereinstimmt.  $\square$

Ziel des restlichen Kapitels ist es, die Optimalität von Tests bei zusammengesetzten Hypothesen zu zeigen. Dies ist vor allem dann möglich, wenn das statistische Modell stochastisch geordnet ist.

**Definition 4.14 (Monotone Dichtequotienten).** Sei  $(\mathcal{P}, \leq)$  total geordnet und  $(X, \{\mathcal{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell und  $T = t(X)$  für  $t : \mathbb{R}^n \rightarrow \mathbb{R}$ . Dann hat  $\mathcal{P}$  einen monotonen Dichtequotienten in  $t$ , wenn für  $\theta, \theta' \in \mathcal{P}$  mit  $\theta \leq \theta'$  eine monotone Abbildung  $p_{\theta, \theta'}$  existiert mit

$$\frac{p_{\theta'}}{p_\theta} = p_{\theta, \theta'} \circ t$$

(( $p_\theta + p_{\theta'}$ )  $\cdot \lambda^n$ -fast sicher).

**Beispiel 4.15 (Exponentialfamilie).** Sei  $p_\theta(x) = h(x) \exp(c(\theta)^\top t(x) - d(\theta))$ , d.h.  $(X, \{\mathbb{P}_\theta = p_\theta \cdot \lambda^n\})$  ist eine ein-parametrische Exponentialfamilie. Dann ist

$$\frac{p_{\theta'}(x)}{p_\theta(x)} = \exp((c(\theta') - c(\theta))^\top t(x) - (d(\theta') - d(\theta))),$$

also hat  $\mathcal{P}$  genau dann einen monotonen Dichtequotienten in  $t$ , wenn  $\mathcal{P}$  total geordnet ist und  $c$  monoton ist.

**Beispiel 4.16 (Beispiel Norm).** Für das Modell aus Beispiel 2.23 mit bekanntem  $\sigma^2 > 0$  ist

$$p_\theta(x) = \exp\left(-\frac{x^2}{2\sigma^2}\right) \cdot \exp\left(\frac{\theta x}{\sigma^2} - \frac{1}{2}\left(\frac{\theta^2}{\sigma^2} + \log(2\pi\sigma^2)\right)\right),$$

es handelt sich also um eine 1-parametrische Exponentialfamilie mit

$$\begin{aligned} c(\theta) &= \frac{\theta}{\sigma^2}, & t(x) &= x, \\ h(x) &= \exp\left(-\frac{x^2}{2\sigma^2}\right), & d(\theta) &= -\frac{1}{2}\left(\frac{\theta^2}{\sigma^2} + \log(2\pi\sigma^2)\right). \end{aligned}$$

insbesondere ist  $c$  monoton, und damit hat  $\mathcal{P}$  einen monotonen Dichtequotienten.

**Proposition 4.17 (Stochastische Monotonie).** Sei  $(\mathcal{P}, \leq)$  total geordnet und  $(X, \{\mathcal{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell und  $T = t(X)$  für  $t : \mathbb{R}^n \rightarrow \mathbb{R}$ . Sei  $f$  isoton (und so, dass  $\mathbb{E}_\theta[f(t(X))]$  für alle  $\theta \in \mathcal{P}$  existiert), und hat  $\mathcal{P}$  einen monotonen Dichtequotienten in  $t$ , dann ist  $\theta \mapsto \mathbb{E}_\theta[f(t(X))]$  ebenfalls isoton.

*Beweis.* Wir schreiben für  $\theta \leq \theta'$

$$\begin{aligned} \mathbb{E}_{\theta'}[f(t(X))] - \mathbb{E}_\theta[f(t(X))] &= \int (f(t(y)) - f(t(x))) p_{\theta'}(y) p_\theta(x) \lambda^n(dx) \lambda^n(dy) \\ &= \int \mathbf{1}_{t(y) > t(x)} ((f(t(y)) - f(t(x))) p_{\theta'}(y) p_\theta(x) \\ &\quad + (f(x) - f(y)) p_{\theta'}(x) p_\theta(y)) \lambda^n(dx) \lambda^n(dy) \\ &= \int \mathbf{1}_{t(y) > t(x)} (f(t(y)) - f(t(x))) (p_{\theta'}(y) p_\theta(x) - p_{\theta'}(x) p_\theta(y)) \lambda^n(dx) \lambda^n(dy) \\ &\geq 0, \end{aligned}$$

da  $f \circ t$  isoton ist und  $\frac{p_{\theta'}(y)}{p_\theta(y)} = p_{\theta, \theta'}(t(y)) \geq p_{\theta, \theta'}(t(x)) = \frac{p_{\theta'}(x)}{p_\theta(x)}$  auf  $\{t(y) > t(x)\}$ .  $\square$

**Theorem 4.18 (Beste Tests bei einseitigen, zusammengesetzten Hypothesen).** Sei  $(\mathcal{P}, \leq)$  total geordnet und  $(X, \{\mathcal{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell,  $T = t(X)$  für  $t : \mathbb{R}^n \rightarrow \mathbb{R}$  und habe  $\mathcal{P}$  einen streng monotonen Dichtequotienten in  $t$ . Weiter sei  $\alpha \in [0, 1]$  und  $\theta_0 \in \mathcal{P}$  so, dass  $\mathcal{P}_0 = \{\theta \leq \theta_0\}$ ,  $\mathcal{P}_1 = \{\theta > \theta_0\}$ . Dann gilt:

1. Sei

$$\psi(x) := \psi_{m, \gamma}(x) := \begin{cases} 1, & t(x) > m, \\ \gamma, & t(x) = m, \\ 0, & t(x) < m \end{cases}$$

für ein  $m \in [0, \infty]$  und  $\gamma \in [0, 1]$ , so dass  $\mathbb{E}_{\theta_0}[\psi(X)] = \alpha$ . Dann ist  $\psi$  bester Test zum Niveau  $\alpha$ .

2. Die Gütefunktion  $\beta_\psi : \theta \mapsto \mathbb{E}_\theta[\psi(X)]$  ist isoton und für  $\theta < \theta_0$  gilt

$$\mathbb{E}_\theta[\psi(X)] = \inf\{\mathbb{E}_\theta[\varphi(X)] : \varphi \in \Phi, \mathbb{E}_{\theta_0}[\varphi(X)] \geq \alpha\},$$

d.h.  $\psi$  minimiert den Fehler erster Art.

*Beweis.* Wir beginnen mit dem Beweis von 1. Sei zunächst  $\theta_1 > \theta_0$ . Wir betrachten den einfachen Test mit  $\mathcal{P}_i = \{\theta_i\}, i = 0, 1$  und zeigen, dass  $\psi$  auch in diesem Fall ein LQ-Test ist. Wegen des streng monotonen Dichtequotienten ist für  $k := p_{\theta_0, \theta_1}(m)$

$$\begin{aligned} \{t(x) > m\} &= \{p_{\theta_0, \theta_1}(t(x)) > p_{\theta_0, \theta_1}(m)\} = \{p_{\theta_1}(x)/p_{\theta_0}(x) > k\}, \\ \{t(x) = m\} &= \{p_{\theta_0, \theta_1}(t(x)) = p_{\theta_0, \theta_1}(m)\} = \{p_{\theta_1}(x)/p_{\theta_0}(x) = k\}, \\ \{t(x) < m\} &= \{p_{\theta_0, \theta_1}(t(x)) < p_{\theta_0, \theta_1}(m)\} = \{p_{\theta_1}(x)/p_{\theta_0}(x) < k\}. \end{aligned}$$

Damit ist  $\psi$  also LQ-Test und  $\mathbb{E}_{\theta_0}[\psi(X)] = \alpha$  nach Voraussetzung. Nach dem Neyman-Pearson-Lemma, Theorem 4.13, ist also  $\psi$  bester Test zum Niveau  $\alpha$ , d.h. es gilt  $\mathbb{E}_{\theta_1}[\psi(X)] = \sup_{\varphi \in \Phi_\alpha} \mathbb{E}_{\theta_1}[\varphi(X)]$ . Da dies aber für alle  $\theta_1 \in \mathcal{P}_1$  gilt, folgt, dass  $\psi$  bester Test zum Niveau  $\alpha$  für  $\mathcal{P}_0 = \{\theta_0\}, \mathcal{P}_1 = \{\theta > \theta_0\}$  ist. Wegen  $\psi = 1_{\{\cdot \geq m\}} \circ t$ , der Monotonie von  $1_{\{\cdot \geq m\}}$  und da  $\mathcal{P}$  einen monotonen Dichtequotienten in  $t$  hat, folgt aus Proposition 4.17, dass  $\theta \mapsto \mathbb{E}_\theta[\psi(X)]$  monoton ist (d.h. die in 2. behauptete Isotonie der Gütefunktion ist gezeigt). Insbesondere gilt für  $\theta \leq \theta_0$ , dass  $\mathbb{E}_\theta[\psi(X)] \leq \mathbb{E}_{\theta_0}[\psi(X)] = \alpha$ , also  $\psi \in \Phi_\alpha$ . Nach Definition ist damit  $\psi$  bester Test zum Niveau für  $\mathcal{P}_0 = \{\theta \leq \theta_0\}, \mathcal{P}_1 = \{\theta > \theta_0\}$ . Es ist noch zu zeigen, dass  $\psi$  den Fehler erster Art minimiert. Sei hierzu  $\theta < \theta_0$ . Genau wie in 1. zeigt man, dass  $1 - \psi$  ein LQ-Test zum Niveau  $1 - \alpha$  ist für  $\mathcal{P}_0 = \{\theta_0\}, \mathcal{P}_1 = \{\theta < \theta_0\}$ . Insbesondere ist  $1 - \psi$  bester Test zum Niveau  $1 - \alpha$ , d.h. für  $1 - \varphi \in \Phi_{1-\alpha}$  (oder  $\mathbb{E}_{\theta_0}[1 - \varphi(X)] \leq 1 - \alpha$ ) ist  $\mathbb{E}_\theta[1 - \varphi] \leq \mathbb{E}_\theta[1 - \psi], \theta < \theta_0$ , also auch  $\mathbb{E}_\theta[\psi] \leq \mathbb{E}_\theta[\varphi]$  für  $\mathbb{E}_\theta[\varphi(X)] \geq \alpha$  und damit die Behauptung.  $\square$

**Korollar 4.19 (Gauß-Test).** Sei  $\mathcal{P} = \mathbb{R}$ ,  $(X, \{\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2)^N\})$  für ein  $\sigma^2 > 0$  wie in (Norm 1b) und  $\theta_0 \in \mathbb{R}$  mit  $\mathcal{P}_0 = (-\infty, \theta_0], \mathcal{P}_1 = (\theta_0, \infty)$ . Weiter sei  $q_\alpha$  das  $\alpha$ -Quantil von  $\mathcal{N}(0, 1)$ . Dann ist

$$\psi(x) := 1\left(\frac{\bar{x} - \theta_0}{\sqrt{\sigma^2/n}} \geq q_\alpha\right)$$

besten Test zum Niveau  $\alpha$ .

*Beweis.* Nach Beispiel 2.23 und Bemerkung 2.28 ist  $\mathcal{N}(\theta, \sigma^2)^n$  eine ein-parametrische Exponentialfamilie mit  $t(x) = \bar{x}$ . Deshalb hat diese Familie nach Proposition 4.17 einen monotonen Dichtequotienten. Weiter ist  $\mathbb{E}_{\theta_0}[\psi(X)] = \mathbb{P}(Z \geq q_\alpha) = \alpha$  für  $Z \sim \mathcal{N}(0, 1)$ . Damit ist  $\psi$  nach Theorem bester Test zum Niveau  $\alpha$ .  $\square$

**Bemerkung 4.20 (Weitere beste Tests).** Wir haben nun den einfachsten Fall eines besten Tests behandelt. Ebenfalls beste Tests sind zweiseitige Gauß-Tests (hier genügen relativ einfache Abschätzungen), aber auch einfache  $t$ -Tests sowie  $\chi^2$ -Tests. Bei den  $t$ -Tests ist  $\sigma^2$  unbekannt, und deshalb ist  $\mathcal{P}$  nicht mehr total geordnet. Deshalb benötigt man dort den Begriff des bedingten Tests, um die Optimalität des Tests zu zeigen.

## 5 Schätztheorie

Im Folgenden sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  mit einem statistischen Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  und  $g : \mathcal{P} \rightarrow \mathfrak{N}$  eine Abbildung auf den Entscheidungsraum  $\mathfrak{N}$  (typischerweise mit  $\mathfrak{N} \subseteq \mathbb{R}^k$  für ein  $k$ ). In diesem Kapitel betrachten wir nur nicht-randomisierte Entscheidungsfunktionen oder Schätzer  $d : E \rightarrow \mathfrak{N}$ . Wir sagen hier, dass  $d$  ein Schätzer für  $g$  ist.

### 5.1 Grundlagen

Es gibt verschiedene Prinzipien, mit denen man zu Schätzern kommt. Wir stellen in diesem Abschnitt das Substitutionsprinzip (und darauf aufbauend die Momentenmethode) und das Maximum-Likelihood-Prinzip vor. Aber erst gibt es ein paar Grundbegriffe.

**Definition 5.1 (Bias, Mean Squared Error).** Sei  $d$  ein Schätzer für  $g$  und  $\mathfrak{N} \subseteq \mathbb{R}$ . Dann heißt

$$b_\theta(d) := \mathbb{E}_\theta[d(X)] - g(\theta)$$

der Bias oder die Verzerrung von  $d$ . Im Fall  $b_\theta(d) = 0$  für alle  $\theta \in \mathcal{P}$  heißt  $d$  unverzerrt oder erwartungstreu oder unbiased. Weiter heißt

$$\mathbb{E}_\theta[(d(X) - g(\theta))^2]$$

die mittlere quadratische Abweichung oder der Mean Squared Error.

**Lemma 5.2 (Zerlegung des Risikos).** Sei  $\mathfrak{N} \subseteq \mathbb{R}$ ,  $d$  ein Schätzer für  $g$  und  $\ell$  der Gauß-Verlust, d.h.  $\ell(\theta, a) = |a - g(\theta)|^2$ . Dann gilt für die Risikofunktion

$$R_d(\theta) = \mathbb{V}_\theta[d(X)] + b_\theta(d)^2.$$

*Beweis.* Wir schreiben

$$\begin{aligned} R_d(\theta) &= \mathbb{E}_\theta[(d(X) - g(\theta))^2] = \mathbb{E}_\theta[(d(X))^2] - \mathbb{E}_\theta[(d(X))^2] + \mathbb{E}_\theta[(d(X))^2] - 2g(\theta)\mathbb{E}_\theta[d(X)] + g(\theta)^2 \\ &= \mathbb{V}_\theta[d(X)] + b_\theta(d)^2. \end{aligned}$$

□

**Definition 5.3 (Plug-in-Schätzer, Maximum-Likelihood-Schätzer).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  ein Schätzproblem.

1. Sei  $X = (X_1, \dots, X_n)$  ein Vektor unabhängiger und identisch verteilter Zufallsvariablen. Dann heißt die (zufällige) Wahrscheinlichkeitsverteilung

$$P_X := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

empirische Verteilung von  $X$ .

2. Sei  $X = (X_1, \dots, X_n)$  ein Vektor identisch verteilter Zufallsvariablen und  $g(\theta) = \mathbb{E}_\theta[f(X_1)]$  für eine Funktion  $f$ , falls der Erwartungswert existiert. Dann heißt

$$\hat{g}(X) := E_X[f] = \frac{1}{n} \sum_{i=1}^n f(X_i)$$

Plug-in-Schätzer für  $g$ . Plug-In-Schätzer für  $g(\theta) = \mathbb{E}_\theta[X_1]$  heißen auch momentenbasierte Schätzer.

3. Sei  $(X, \{\mathbb{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  ein reguläres statistisches Problem und  $\aleph = \mathcal{P}$ . Existiert eine messbare Abbildung  $\hat{\theta} : E \rightarrow \aleph$  mit

$$p_{\hat{\theta}(x)}(x) = \sup_{\theta \in \mathcal{P}} p_\theta(x),$$

so heißt  $\hat{\theta}$  Maximum-Likelihood-Schätzer von  $g(\theta) = \theta$ . Weiter heißt für festes  $x$  die Funktion  $\theta \mapsto p_\theta(x)$  Likelihood und  $\theta \mapsto \log p_\theta(x)$  Log-Likelihood.

**Bemerkung 5.4 (Einfache Eigenschaften).** 1. Sei  $X = (X_1, \dots, X_n)$  ein Vektor identisch verteilter Zufallsvariablen. Dann ist jeder Plug-In-Schätzer  $\hat{g}$  für  $g$  (mit  $g(\theta) := \mathbb{E}_\theta[f(X_1)]$ ) unverzerrt.

Denn: Es gilt

$$\mathbb{E}_\theta[\hat{g}(X)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\theta[f(X_i)] = \mathbb{E}_\theta[f(X_1)] = g(\theta).$$

2. Sei  $d : E \rightarrow \aleph$  ein Schätzer für  $g : \mathcal{P} \rightarrow \aleph$  und  $f : \aleph \rightarrow \aleph$ . Ist  $d$  ein Maximum-Likelihood-Schätzer für  $g$ , so ist  $d \circ f$  ein Maximum-Likelihood-Schätzer für  $g \circ f$ . Ist  $d$  unverzerrt, so ist jedoch der Schätzer  $d \circ f$  für  $g \circ f$  nicht unbedingt unverzerrt.

Denn: Der Wert, an dem eine Funktion ihr Maximum annimmt, verändert sich nicht, wenn man die Funktion mit einer weiteren Funktion verknüpft, was die erste Behauptung zeigt. Für die zweite Behauptung sei etwa  $d$  ein Plug-In-Schätzer für  $g(\theta) = \mathbb{E}_\theta[X_1]$ , sowie  $f : x \mapsto x^2$ . Wäre  $d \circ f$  unverzerrt, so müsste

$$\mathbb{E}_\theta[d(X)^2] = \mathbb{E}_\theta[f(d(X))] = f(g(\theta)) = (g(\theta))^2 = (\mathbb{E}_\theta[d(X)])^2,$$

also  $\mathbb{V}_\theta[d(X)] = 0$ , was im Allgemeinen sicher falsch ist.

**Proposition 5.5 (Maximum-Likelihood-Schätzer in Exponentialfamilien).** Sei  $(X, \{\mathbb{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  eine  $k$ -parametrische Exponentialfamilie mit  $c, t, d, h$  für  $c_1, \dots, c_k$  injektiv, also

$$p_\theta(x) = h(x) \cdot \exp(c(\theta)^\top t(x) - d(\theta)).$$

Sei  $C = \{c(\theta) : \theta \in \mathcal{P}\}^\circ$ . Falls die Gleichung (in  $\theta$ )

$$\mathbb{E}_\theta[t_i(X)] = t_i(x), \quad i = 1, \dots, k$$

(in  $\theta$ ) eine Lösung  $x \mapsto \hat{\theta}(x)$  besitzt mit  $c(\hat{\theta}(x)) \in C$ , dann ist  $\hat{\theta}$  der eindeutige Maximum-Likelihood-Schätzer für  $\theta$ .

*Beweis.* Wir zeigen die Behauptung nur für  $k = 1$ . Sei oBdA  $c = \text{id}$ . Den allgemeinen Fall zeigt man mit Bemerkung 5.4.2. Wir berechnen

$$\frac{\partial}{\partial \theta_i} \log p_\theta(x) = t(x) - d'(\theta), \quad \frac{\partial^2}{\partial \theta_i^2} \log p_\theta(x) = -d''(\theta).$$

Daraus folgt mit Proposition 2.35, dass  $\mathbb{E}_\theta[t(X)] = d'(\theta)$  und  $\mathbb{V}_\theta[t(X)] = -d''(\theta)$ . Damit ist gezeigt, dass die Log-Likelihood-Funktion wegen der negativen Ableitung strikt konkav ist und  $\hat{\theta}$  genau dann ein Maximum-Likelihood-Schätzer ist, wenn  $t(x) = \mathbb{E}_{\hat{\theta}(x)}[t(X)]$  gilt. (Die Eindeutigkeit folgt mit der Konkavität.)  $\square$

**Beispiel 5.6 (Beispiel Norm).** Für das Beispiel aus Beispiel 2.23 erinnern wir an Bemerkung 2.28 und sehen, dass es sich um eine 2-parametrische Exponentialfamilie handelt mit  $t_1(x) = \sum_{i=1}^n x_i$ ,  $t_2(x) = \sum_{i=1}^n x_i^2$ . Wir betrachten also etwa die Gleichung

$$\mathbb{E}_{(\mu(x), \sigma^2(x))}[t_1(X)] = t_1(x),$$

die genau dann erfüllt ist, wenn

$$n\mu(x) = t_1(x), \quad \text{also für } \mu(x) = \bar{x}.$$

Damit ist bereits  $\bar{X}$  der maximum-Likelihood-Schätzer für  $\mu$ . Weiter betrachten wir die Gleichung

$$\mathbb{E}_{(\mu(x), \sigma^2(x))}[t_2(X)] = t_2(x)$$

die genau dann erfüllt ist, wenn

$$n(\sigma^2(x) + \mu^2(x)) = t_2(x), \quad \text{also für } \sigma^2(x) = \frac{1}{n}t_2(x) - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Damit haben wir den maximum-Likelihood-Schätzer für  $\sigma^2$  hergeleitet.

**Beispiel 5.7 (Lineare Regression).** Wir kehren zurück zur linearen Regression aus Beispiel 2.30. Wieder leiten wir Maximum-Likelihood-Schätzer für die Parameter  $\beta_0, \dots, \beta_m, \sigma^2$  her. Für  $t(y)^\top = (t_0(y), \dots, t_m(y))$  schreiben wir

$$\mathbb{E}_{(\beta(y), \sigma^2(y))}[t(Y)]^\top = xx^\top \beta(y) = t(y)^\top = xy,$$

$$\mathbb{E}_{(\beta(y), \sigma^2(y))}[t_{m+1}(Y)] = (x^\top \beta)^2 + \sigma^2(y) = t_{m+1}(y) = \sum_{i=1}^n y_i^2 = yy^\top.$$

Aus der ersten Gleichung lesen wir ab, dass

$$\hat{\beta}(y) = (xx^\top)^{-1}xy$$

und aus der zweiten, dass

$$\widehat{\sigma^2}(y) = yy^\top - (x^\top (xx^\top)^{-1}xy)^2$$

die Maximum-Likelihood-Schätzer für die Parameter  $\beta$  und  $\sigma^2$  sind.

## 5.2 UMVUE-Schätzer

Um es vorwegzunehmen: UMVUE steht für *Uniformly Minimum Variance Unbiased Estimator*. Solche Schätzer minimieren also unter allen unverzerrten Schätzern die Varianz.

**Definition 5.8 (UMVUE-Schätzer).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  ein Schätzproblem. Ein Schätzer  $d$  für  $g$  heißt UMVUE oder Uniformly Minimum Variance Unbiased Estimator für  $g$ , falls er unverzerrt ist und

$$\mathbb{V}_\theta[d(X)] = \inf_{e \text{ unverzerrt}} \mathbb{V}_\theta[e(X)], \quad \theta \in \mathcal{P}.$$



**Bemerkung 5.9 (Rao-Blackwell).** Wir erinnern kurz an den Satz von Rao-Blackwell, Theorem 3.10, für ein Schätzproblem  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  im Falle des Gauß-Verlustes. Sei  $T = t(X)$  suffizient und  $d$  ein nicht-randomisierter Schätzer für  $g$  mit  $\mathbb{E}_\theta[|d(X)|] < \infty$  für alle  $\theta \in \mathcal{P}$ . Dann hat der Schätzer  $e \circ t$  für  $g$  mit

$$e(t) := \mathbb{E}_\theta[d(X)|T = t]$$

geringeres Risiko als  $d$ .

**Theorem 5.10 (Lehmann-Scheffé).** Sei  $((X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  ein Schätzproblem mit Gauß-Verlust  $\ell$ ,  $T = t(X)$  eine vollständige, suffiziente Statistik und  $d$  ein unverzerrter Schätzer für  $g$ . Dann ist  $e \circ t$  mit

$$e(t) := \mathbb{E}_\theta[d(X)|T = t]$$

ein UMVUE für  $g$ . Ist außerdem  $\mathbb{V}_\theta[e(t((X)))] < \infty$  für alle  $\theta \in \mathcal{P}$ , so ist  $e \circ t$  der einzige UMVUE für  $g$ .

*Beweis.* Da  $d$  unverzerrt ist, folgt mit der Turmeigenschaft der bedingten Erwartung, dass auch  $e \circ t$  unverzerrt ist. Wir zeigen nun, dass  $e$  und damit  $e \circ t$  unabhängig von der Wahl von  $d$  ist. Dann hat insbesondere  $e \circ t$  minimale Varianz unter allen unverzerrten Schätzern, ist also UMVUE. Seien also  $d_1, d_2$  unverzerrte Schätzer für  $g$ . Weiter seien  $e_i \circ t$  für  $e_i(t) = \mathbb{E}_\theta[d_i(X)|T = t]$  zwei unverzerrte Schätzer. Insbesondere gilt

$$\mathbb{E}_\theta[e_1(T) - e_2(T)] = \mathbb{E}_\theta[d_1(X) - d_2(X)] = g(\theta) - g(\theta) = 0$$

für alle  $\theta \in \mathcal{P}$ . Da  $T$  vollständig ist folgt  $e_1(T) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} e_2(T)$  für alle  $\theta \in \mathcal{P}$ , was zu zeigen war. Für die Eindeutigkeit sei  $d'$  ein weiterer UMVUE. Einerseits ist dann nach dem Satz von Rao-Blackwell

$$\mathbb{E}_\theta[(\mathbb{E}_\theta[d'(X)|T] - g(\theta))^2] \leq \mathbb{E}_\theta[(d'(X) - g(\theta))^2],$$

und da  $d'$  ein UMVUE ist, gilt hier Gleichheit, also  $d'(X) = \mathbb{E}_\theta[d'(X)|T]$ . Andererseits haben wir oben gezeigt, dass  $e(T) = \mathbb{E}_\theta[d'(X)|T]$ , da  $d'$  unverzerrt ist. Insgesamt gilt also  $d' = e \circ t$ .  $\square$

**Beispiel 5.11 (Beispiel Unif).** Wir betrachten das Schätzproblem mit  $g(\theta) = \theta$  und Gauß-Verlust  $\ell$ . Hierzu verwendet wir die Statistik  $T = t(X) = \max_{1 \leq i \leq n} X_i$ . Wir haben bereits in Beispiel 2.8 und 2.17 gesehen, dass  $T$  vollständig und suffizient ist. Aus Beispiel 3.11 wissen wir außerdem, dass

$$d(X) := \frac{n+1}{n} \max_{1 \leq i \leq n} X_i$$

unverzerrt ist und dieser Schätzer eine kleinere Risikofunktion hat als  $2\bar{x}$ . Außerdem gilt  $\mathbb{E}_\theta[d(X)|T] = d(X)$ , da  $d(X)$  eine Funktion von  $T$  ist. Nach dem Satz von Lehmann und Scheffé ist also  $d$  ein UMVUE. Wegen  $\mathbb{V}_\theta[d(X)] < \infty$  für alle  $\theta \in \mathcal{P}$  ist dieser UMVUE auch einseitig.

Die Argumentation des letzten Beispiels lässt sich verallgemeinern, was wir nun für Exponentialfamilien tun wollen.

**Korollar 5.12 (UMVUE bei Exponentialfamilien).** Sei  $(X, \{\mathbb{P}_\theta = p_\theta \cdot \lambda^n : \theta \in \mathcal{P}\})$  eine  $k$ -parametrische Exponentialfamilie mit  $c, t, d, h$ , also

$$p_\theta(x) = h(x) \cdot \exp(c(\theta)^\top t(x) - d(\theta)).$$

Weiter sei  $g(\theta) = \mathbb{E}_\theta[f(t(X))]$ . Dann ist  $d(X) := f(t(X))$  ein UMVUE für  $g$ . Ist außerdem  $\mathbb{V}_\theta[d(X)] < \infty$  für alle  $\theta \in \mathcal{P}$ , so gibt es nur diesen einen UMVUE.

*Beweis.* Zunächst wissen wir aus Proposition 2.29 und Theorem 2.32, dass  $T = t(X) = (t_1(X), \dots, t_k(X))$  suffizient und vollständig ist. Weiter ist  $d$  unverzerrt und  $\mathbb{E}_\theta[d(X)|T] = d(X)$ , da  $d(X)$  messbar bezüglich  $T$  ist. Nach dem Satz von Lehmann und Scheffé ist damit  $d$  ein UMVUE für  $g$ . Die letzte Behauptung wurde schon in Theorem 5.10 bewiesen.  $\square$

**Beispiel 5.13 (Beispiel Norm).** Für das Beispiel *Norm* aus Beispiel 2.23 haben wir in Beispiel 5.6 gesehen, dass es sich um eine 2-parametrische Exponentialfamilie handelt mit  $t_1(x) = \sum_{i=1}^n x_i$ ,  $t_2(x) = \sum_{i=1}^n x_i^2$ . Wir betrachten nun  $g_1(\mu, \sigma^2) = \mu$  und  $g_2(\mu, \sigma^2) = \sigma^2$  und schreiben

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \left( \sum_{i=1}^n x_i^2 \right) - 2n\bar{x}^2 + n\bar{x}^2 \right) = \frac{1}{n-1} t_2(x) - \frac{1}{(n-1)n} t_1(x)^2,$$

also

$$\mathbb{E}_\theta \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \frac{n}{n-1} \sigma^2 - \frac{1}{n-1} \sigma^2 = \sigma^2$$

und damit sind

$$\begin{aligned} g_1(\theta) &= \mathbb{E}_\theta \left[ \frac{1}{n} t_1(X) \right] = \mathbb{E}_\theta[\bar{X}] = \mu, \\ g_2(\theta) &= \mathbb{E}_\theta \left[ \frac{1}{n-1} t_2(X) - \frac{1}{(n-1)n} t_1^2(X) \right] = \sigma^2 \end{aligned}$$

eindeutige UMVUEs.

**Beispiel 5.14 (Lineare Regression).** Wir betrachten noch einmal die lineare Regression aus Beispiel 2.30. Für  $t(y)^\top = (t_0(y), \dots, t_m(y))$  schreiben wir für  $f(t(y)) = (xx^\top)^{-1}t(y)$

$$\mathbb{E}_{(\beta, \sigma^2)}[f(t(Y))] = \mathbb{E}_{(\beta, \sigma^2)}[(xx^\top)^{-1}t(Y)] \mathbb{E}_{(\beta, \sigma^2)}[(xx^\top)^{-1}xY] = (xx^\top)^{-1}xx^\top \beta = \beta.$$

Damit ist  $(xx^\top)^{-1}xy$  ein UMVUE für  $\beta$  (und ein Maximum-Likelihood-Schätzer nach Beispiel 5.7).

### 5.3 Information und die Cramér-Rao-Schranke

Zwar haben wir bereits obere Schranken für das Risiko eines Schätzers ermittelt, jedoch gibt es auch untere Schranken, insbesondere die in Theorem 5.22 vorgestellte Cramér-Rao-Schranke. Für diese benötigen wir den Begriff der (Fisher-)Information. Diese beschreibt die Krümmung (im Sinne der zweiten Ableitung) der log-Likelihood in Bezug auf die Parameter. Ist diese Zahl groß, bedeutet das eine hohe Krümmung, und damit ist die log-Likelihood für nahe Parameter deutlich kleiner. Anders ausgedrückt haben wir über die log-Likelihood viel über den wahren Parameter gelernt, es gibt also viel "Information". Um die zweite Ableitung sinnvoll berechnen zu können, benötigen wir zunächst ein paar Regularitätsannahmen.

**Annahme 5.15 (Regularitätsannahmen der Fisher-Information).** Für ein reguläres statistisches Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  auf  $\mathbb{R}^n$  mit  $\mathcal{P} \subseteq \mathbb{R}^m$  und  $\mathbb{P}_\theta = p_\theta \cdot \lambda^n$  geben wir folgende Regularitätsannahmen an:

(A1) Es gibt ein  $N \in \mathcal{B}(\mathbb{R}^n)$ , so dass für alle  $\theta \in \mathcal{P}$  und  $i = 1, \dots, k$  die Ableitung  $\partial p_\theta(x)/\partial \theta_i$  für  $x \notin N$  existiert.

(A2) Für jedes messbare  $\phi$  gilt, falls der Erwartungswert existiert,

$$\frac{\partial}{\partial \theta_i} \mathbb{E}_\theta[\phi(X)] = \int \phi(x) \frac{\partial p_\theta(x)}{\partial \theta_i} d\lambda^n(x).$$

(A3) Die Menge

$$C := \{x : p_\theta(x) > 0\}$$

ist unabhängig von  $\theta$ .

(A4) Es gilt (A2) und außerdem auch

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \mathbb{E}_\theta[\phi(X)] = \int \phi(x) \frac{\partial^2 p_\theta(x)}{\partial \theta_i \partial \theta_j} d\lambda^n(x)$$

für jedes messbare  $\phi$ , für das der Erwartungswert existiert.

**Beispiel 5.16 (Beispiele *Bern*, *Norm*, *Unif*).** Für Exponentialfamilien ist Annahme 5.15 nach Lemma 2.34 immer erfüllt, insbesondere also für Beispiele *Bern* und *Norm*. Für Beispiel 1.8 ist jedoch sowohl (A1) als auch (A3) verletzt.

**Definition 5.17 (Fisher-Information).** Sei  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  ein statistisches Modell, das die Annahmen (A1)–(A3) erfüllt. Dann heißt

$$I(\theta) := I(\mathbb{P}_\theta) := \left( \mathbb{C}\text{OV} \left[ \frac{\partial}{\partial \theta_i} \log p_\theta(X), \frac{\partial}{\partial \theta_j} \log p_\theta(X) \right] \right)_{i,j=1,\dots,k}$$

Fisher-Informations-Matrix.

**Bemerkung 5.18 (Berechnung der Fisher-Information unter (A4)).** Gilt statt (A2) sogar (A4), so kann man schreiben

$$\begin{aligned} \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p_\theta(X) \right] &= \mathbb{E}_\theta \left[ \frac{1}{p_\theta(X)^2} \left( p_\theta(X) \frac{\partial^2 p_\theta(X)}{\partial \theta_i \partial \theta_j} - \frac{\partial p_\theta(X)}{\partial \theta_i} \cdot \frac{\partial p_\theta(X)}{\partial \theta_j} \right) \right] \\ &= 0 - \mathbb{E}_\theta \left[ \frac{\partial \log p_\theta(X)}{\partial \theta_i} \frac{\partial \log p_\theta(X)}{\partial \theta_j} \right] = -(I(\theta))_{i,j}. \end{aligned}$$

Dies liefert also eine alternative Berechnung der Fisher-Information.

**Beispiel 5.19 (Exponentialfamilie).** Im Falle einer Exponentialfamilie in kanonischer Form ist

$$p_\theta(x) = h(x) \exp(\theta^\top t(x) - d(\theta))$$

und damit mit Proposition 2.35

$$\frac{\partial}{\partial \theta_i} \log p_\theta(x) = t_i(x) - \frac{\partial d(\theta)}{\partial \theta_i} = t_i(X) - \mathbb{E}_\theta[t_i(X)].$$

Weiter ist

$$(I(\theta))_{i,j} = \mathbb{C}\text{OV}_\theta[t_i(X), t_j(X)].$$

**Beispiel 5.20 (Fisher-Information einer unabhängigen Stichprobe).** Für ein statistisches Modell  $(X, \{\mathbb{P}_\theta : \theta \in \mathcal{P}\})$  haben wir die Fisher-Information definiert. Nun betrachten wir den Fall von Daten  $X_1, \dots, X_n$ , bei denen die  $X_i$  unter allen  $\mathbb{P}_\theta$  unabhängig sind. Anders gesagt betrachten wir das statistische Modell  $(X = (X_1, \dots, X_n), \{\mathbb{P}_\theta^n : \theta \in \mathcal{P}\})$ , wobei  $\mathbb{P}_\theta^n$  das  $n$ -fache Produktmaß ist. Hier ist die Dichte gegeben durch  $x = (x_1, \dots, x_n) \mapsto p_\theta(x_1) \cdots p_\theta(x_n)$  und damit gilt

$$\begin{aligned} I(\mathbb{P}_\theta^n) &= \text{COV}_\theta^n \left[ \frac{\partial}{\partial \theta_j} \log p_\theta(X), \frac{\partial}{\partial \theta_k} \log p_\theta(X) \right]_{jk} \\ &= \text{COV}_\theta^n \left[ \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \log p_\theta(X_i), \sum_{i=1}^n \frac{\partial}{\partial \theta_k} \log p_\theta(X_i) \right]_{jk} \\ &= \sum_{i=1}^n \text{COV}_\theta \left[ \frac{\partial}{\partial \theta_j} \log p_\theta(X_i), \frac{\partial}{\partial \theta_k} \log p_\theta(X_i) \right]_{jk} \\ &= nI(\mathbb{P}_\theta). \end{aligned}$$

**Proposition 5.21 (Mittlere Ableitung der log-Likelihood verschwindet).** *Es gelte (A1)–(A3). Dann ist für  $i = 1, \dots, m$*

$$\mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta_i} \log p_\theta(X) \right] = 0.$$

*Beweis.* Wir schreiben

$$\begin{aligned} \mathbb{E}_\theta \left[ \frac{\partial}{\partial \theta_i} \log p_\theta(X) \right] &= \int p_\theta(x) \frac{\partial}{\partial \theta_i} \log p_\theta(x) d\lambda^n(dx) = \int p_\theta(x) \frac{1}{p_\theta(x)} \frac{\partial}{\partial \theta_i} p_\theta(x) d\lambda^n(dx) \\ &= \frac{\partial}{\partial \theta_i} \mathbb{E}_\theta[1] = 0, \end{aligned}$$

wobei die Vertauschung von Integral und Ableitung nach (A3) erlaubt ist.  $\square$

**Theorem 5.22 (Cramér-Rao-Schranke).** *Es gelte (A1)–(A3) sowie  $\mathcal{P} \subseteq \mathbb{R}$  und  $I(\theta) > 0$  für alle  $\theta \in \mathcal{P}$ . Sei  $T = t(X)$  für  $t : \mathbb{R}^m \rightarrow \mathbb{R}$  eine Statistik und es existiere  $\Psi(\theta) := \mathbb{E}_\theta[T]$ . Dann gilt*

$$\mathbb{V}_\theta[T] \geq \frac{\Psi'(\theta)^2}{I(\theta)}.$$

*Ist insbesondere  $\mathbb{E}_\theta[T] = \theta$  (d.h.  $T$  ist ein erwartungstreuer Schätzer für  $\theta$ ), so gilt*

$$\mathbb{V}_\theta[T] \geq \frac{1}{I(\theta)}.$$

*Beweis.* Wir schreiben mit Proposition 5.21 und der Cauchy-Schwartz-Ungleichung

$$\begin{aligned} \Psi'(\theta) &= \int t(x) \frac{\partial}{\partial \theta} p_\theta(x) \lambda^n(dx) = \mathbb{E}_\theta \left[ t(x) \frac{\partial}{\partial \theta} \log p_\theta(x) \right] \\ &= \mathbb{E}_\theta \left[ (t(X) - \mathbb{E}_\theta[T]) \frac{\partial}{\partial \theta} \log p_\theta(X) \right] \\ &\leq (\mathbb{V}_\theta[T] \cdot I(\theta))^{1/2}. \end{aligned}$$

Daraus folgt die Behauptung.  $\square$

**Beispiel 5.23 (1-parametrische Exponentialfamilie).** Wir betrachten den Fall aus Beispiel 5.19 mit  $\theta \in \mathbb{R}$  und der suffizienten Statistik  $T = t(X)$ . Wir wissen (etwa aus Beispiel 5.19 und Proposition 2.35), dass  $\mathbb{E}_\theta[t(X)] = d'(\theta)$  und  $I(\theta) = \mathbb{V}_\theta[t(X)] = d''(\theta)$ . Damit gilt

$$\mathbb{V}_\theta[t(X)] = d''(\theta) = \frac{d''(\theta)^2}{I(\theta)} = \frac{(\frac{d}{d\theta}\mathbb{E}_\theta[t(X)])^2}{I(\theta)},$$

und damit gilt Gleichheit in der Schranke von Theorem 5.22.

Andersherum ist im Beweis der Cramér-Rao-Schranke klar: Gleichheit gilt genau dann, wenn Gleichheit in der Cauchy-Schwartz-Ungleichung gilt, also wenn  $t(X) - \mathbb{E}_\theta[T]$  und  $\frac{\partial}{\partial\theta} \log p_\theta(X)$  linear abhängig sind, also wenn für geeignete  $a$  und  $b$

$$\frac{\partial}{\partial\theta} \log p_\theta(x) = a(\theta)t(x) + b(\theta)$$

oder

$$\log p_\theta(x) = t(x) \int_{\theta_0}^{\theta} a(\eta) d\eta + \int_{\theta_0}^{\theta} b(\eta) d\eta.$$

Das bedeutet, dass  $p_\theta$  die Dichte einer Exponentialfamilie ist. Die Cramér-Rao-Schranke ist damit (falls (A1)–(A3) gelten) genau für 1-parametrische Exponentialfamilien scharf.

**Beispiel 5.24 (Beispiel *Unif*).** Aus Beispiel 2.25 wissen wir, dass es sich nicht um eine Exponentialfamilie handelt. Außerdem sind die Regularitätsannahmen (A1) und (A3) nicht erfüllt; siehe Beispiel 5.16. Deshalb ist nicht klar, ob die Cramér-Rao-Schranke in diesem Modell existiert. Um zu sehen, ob die Schranke doch funktioniert, verwenden wir Bemerkung 5.18, so dass

$$I(\theta) = -\frac{\partial^2}{\partial\theta^2} \mathbb{E}_\theta[\log p_\theta(X)] = 1/\theta^2. \quad (5.1)$$

Wir setzen  $t(X) = X$ , so dass  $\mathbb{V}_\theta[t(X)] = \theta^2/12$ . Außerdem ist  $\frac{d}{d\theta} \mathbb{E}_\theta[X] = 1/2$  und damit

$$\mathbb{V}_\theta[t(X)] = \frac{\theta^2}{12} < \frac{\theta^2}{4} = \frac{(\frac{d}{d\theta} \mathbb{E}_\theta[X])^2}{I(\theta)}.$$

Also gilt die Cramér-Rao-Schranke hier nicht.

**Bemerkung 5.25 (Unabhängige Stichprobe).** Für das statistische Modell  $(X^n = (X_1, \dots, X_n), \{\mathbb{P}_\theta^n : \theta \in \mathcal{P}\})$  der  $n$ -fachen unabhängigen Versuchswiederholung,  $n = 1, 2, \dots$  aus Beispiel 5.20 erhalten wir folgendes Resultat (falls (A1)–(A3) gelten):

Ist  $\mathcal{P} \subseteq \mathbb{R}$  und ist  $d^n$  ein Schätzer für  $\theta$  im  $n$ -ten Modell, so gilt

$$\mathbb{V}_\theta[d^n(X^n)] \geq \frac{(\frac{d}{d\theta} \mathbb{E}_\theta[d^n(X^n)])^2}{nI(\theta)}.$$

Asymptotisch bedeutet das, falls  $\sqrt{n}(\mathbb{E}_\theta[d^n(X^n)] - g(\theta)) \xrightarrow{n \rightarrow \infty} 0$  und  $\frac{d}{d\theta}(\mathbb{E}_\theta[d^n(X^n)] - g(\theta)) \xrightarrow{n \rightarrow \infty} 0$

$$\begin{aligned} \liminf_{n \rightarrow \infty} n \mathbb{E}_\theta[(d^n(X^n) - g(\theta))^2] &= \liminf_{n \rightarrow \infty} n \mathbb{V}_\theta[(d^n(X^n)] + n(\mathbb{E}_\theta[d^n(X^n)] - g(\theta))^2 \\ &= \liminf_{n \rightarrow \infty} n \mathbb{V}_\theta[(d^n(X^n)] \geq \liminf_{n \rightarrow \infty} \frac{(\frac{d}{d\theta} \mathbb{E}_\theta[d^n(X^n)])^2}{I(\theta)} = \frac{(g'(\theta))^2}{I(\theta)}. \end{aligned}$$

## 5.4 Asymptotik von Maximum-Likelihood-Schätzern

Das letzte Resultat behandelt schon die Asymptotik einer Folge von Schätzern. Dies wollen wir nun noch für Maximum-Likelihood-Schätzer ausbauen. Einfach gesagt konvergieren Maximum-Likelihood-Schätzer (unter einigen Regularitätsannahmen) für unabhängige Stichproben gegen den wahren Parameter (Theorem 5.30). Noch genauer sind sie asymptotisch normalverteilt und die Inverse der Fisher-Information gibt die Covarianzmatrix (Theorem 5.31). Es folgen nun weitere Annahmen, insbesondere die der unabhängigen, identisch verteilten Daten.

**Bemerkung 5.26 (Weitere Annahmen).** (A5) Es ist  $(X^n = (X_1^n, \dots, X_n^n), \{\mathbb{P}_\theta^n : \theta \in \mathcal{P}\})$  ein reguläres statistisches Modell, so dass  $X^n$  ein unabhängiger, identisch verteilter Vektor ist mit

$$\mathbb{P}^n(X_1^n \in \cdot) = p_\theta \cdot \lambda^m.$$

Insbesondere hat  $\mathbb{P}^n$  die Dichte

$$\prod_{i=1}^n p_\theta(x_i).$$

(A6) Für alle  $x$  ist

$$\eta \mapsto \frac{\partial^2}{\partial \eta_i \partial \eta_j} \log p_\eta(x)$$

stetig und endlich und es gilt für alle  $\theta \in \mathcal{P}$

$$\eta \mapsto \frac{\partial^2}{\partial \eta_i \partial \eta_j} \mathbb{E}_\theta \left[ \log p_\eta(X) \right] = \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \eta_i \partial \eta_j} \log p_\eta(X) \right].$$

**Bemerkung 5.27 (Maximum-Likelihood-Schätzer).** Gilt (A1)–(A5), so ist der Maximum-Likelihood-Schätzer  $d^n(x^n)$  für  $g$ , falls er im Inneren von  $\mathcal{P}$  liegt, Lösung der Gleichung

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \log p_\theta(x_i^n) \Big|_{\theta=d^n(x^n)} = 0. \quad (\text{Lik})$$

**Definition 5.28 (Konsistenz).** Sei  $((X^n, \{\mathbb{P}_\theta^n : \theta \in \mathcal{P}\}), \mathfrak{N}, g, \ell)$  ein Schätzproblem, und  $d^n$  ein Schätzer für  $g$ ,  $n = 1, 2, \dots$ . Dann heißt die Folge  $d^n$  konsistent, falls

$$\mathbb{P}_\theta(|d^n(X^n) - g(\theta)| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$$

für alle  $\varepsilon > 0$  und alle  $\theta \in \mathcal{P}$ . Dies ist insbesondere dann der Fall, wenn  $X^1, X^2, \dots$  auf demselben Wahrscheinlichkeitsraum definiert sind und  $d^n(X^n) \xrightarrow{n \rightarrow \infty}_{f_s} g(\theta)$ .

**Bemerkung 5.29 (Konsistenz bei unabhängigen Daten).** Unter Annahme (A5) ist unter  $\mathbb{P}_\theta^n$  der Vektor  $X^n$  unabhängig und identisch nach  $\mathbb{P}_\theta$  verteilt. Wir können dann oBdA annehmen, dass alle  $X^n$  auf demselben Wahrscheinlichkeitsraum definiert sind. Wir schreiben dann auch  $X$  anstelle von  $X^n$  (und denken uns, dass  $d^n$  ja nur von den ersten  $n$  Einträgen von  $X$  abhängt). Insbesondere kann dann die fast sichere Konvergenz  $d^n(X) \xrightarrow{n \rightarrow \infty} g(\theta)$  unter  $\mathbb{P}_\theta$  gelten.

**Theorem 5.30 (Konsistenz von Maximum-Likelihood-Schätzern).** Sei  $\theta \in \mathcal{P}$  und  $d^n(X)$  der Maximum-Likelihood-Schätzer für  $\theta$  (basierend auf den ersten  $n$  Beobachtungen). Sei

$$Z(M, x) := \inf_{\eta \in M} \log \frac{p_\theta(x)}{p_\eta(x)}.$$

Angenommen, für jedes  $\eta \neq \theta$  gibt es ein  $\varepsilon(\eta) > 0$ , so dass  $\mathbb{E}_\theta[Z(B_{\varepsilon(\eta)}, X)] > 0$ . Weiter existiere ein  $\mathcal{C}$  kompakt mit  $\theta \in \mathcal{C}$  und  $\mathbb{E}_\theta[Z(\mathcal{P} \setminus \mathcal{C}, X)] > 0$ . Dann ist

$$\lim_{n \rightarrow \infty} d^n(X) \stackrel{\mathbb{P}_\theta\text{-fs}}{=} \theta.$$

*Beweis.* Für beliebiges  $\varepsilon > 0$  müssen wir zeigen, dass

$$\mathbb{P}_\theta(\limsup_{n \rightarrow \infty} |d^n(X) - \theta| > \varepsilon) = 0.$$

(Die Aussage folgt dann mit einem Grenzwert  $\varepsilon \downarrow 0$ .) Da  $\mathcal{C} \setminus B_\varepsilon(\theta)$  kompakt ist und  $\{B_{\varepsilon(\eta)}(\eta) : \eta \in \mathcal{C} \setminus B_\varepsilon(\theta)\}$  eine offene Überdeckung der kompakten Menge  $\mathcal{C} \setminus B_\varepsilon(\theta)$  ist, gibt es eine endliche Teilüberdeckung  $A_1 := B_{\varepsilon(\eta_1)}(\eta_1), \dots, A_m := B_{\varepsilon(\eta_m)}(\eta_m)$ . Mit  $A_0 := \mathcal{P} \setminus \mathcal{C}$  ist  $\mathcal{P} = B_\varepsilon(\theta) \cup A_0 \cup \dots \cup A_m$  und  $\mathbb{E}_\theta[Z(A_j, X_i)] =: c_j > 0$ ,  $j = 0, \dots, m$ . Nach dem starken Gesetz der großen Zahlen ist

$$\frac{1}{n} \sum_{i=1}^n Z(A_j, X_i) \xrightarrow{n \rightarrow \infty, \mathbb{P}_\theta\text{-fs}} c_j.$$

Dann gilt

$$\begin{aligned} & \mathbb{P}_\theta(\limsup_{n \rightarrow \infty} |d^n(X) - \theta| \geq \varepsilon) \\ & \leq \mathbb{P}_\theta\left(\text{Es gibt } \eta_1, \eta_2, \dots \notin B_\varepsilon(\theta), \text{ so dass} \right. \\ & \quad \left. \frac{1}{n} \sum_{i=1}^n \log p_{\eta_i}(X_i) \geq \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i) \text{ für unendlich viele } n\right) \\ & \leq \sum_{j=0}^m \mathbb{P}\left(\inf_{\eta \in A_j} \frac{1}{n} \sum_{i=1}^n \log \frac{p_\theta(X_i)}{p_\eta(X_i)} \leq 0 \text{ unendlich oft}\right) \\ & \leq \sum_{j=0}^m \mathbb{P}_\theta\left(\frac{1}{n} \sum_{i=1}^n Z(A_j, X_i) \leq 0 \text{ unendlich oft}\right) = 0 \end{aligned}$$

und die Behauptung ist gezeigt.  $\square$

**Theorem 5.31 (Asymptotische Normalität des Maximum-Likelihood-Schätzers).** Es gelte (A1)–(A6) und es sei für jedes  $\theta \in \mathcal{P}^{\circ 10}$

$$\sup_{j,k} \mathbb{E}_\theta \left[ \sup_{|\eta - \theta| < r} \left| \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X_1) \Big|_{\eta=\theta} - \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X_1) \right| \right] \xrightarrow{r \rightarrow 0} 0. \quad (*)$$

Weiter sei  $d^n(X)$  konsistent und die Fisher-Information  $I(\theta)$  sei für alle  $\theta \in \mathcal{P}$  nicht-singulär. Dann gilt für  $X \sim \mathbb{P}_\theta^\infty$

$$\sqrt{n}(d^n(X) - \theta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, I^{-1}(\theta)).$$

<sup>10</sup>Für eine Menge  $A$  sei  $A^\circ$  das Innere.

*Beweis.* Sei

$$\ell_x^n(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i).$$

Dann ist, wegen der stetigen Differenzierbarkeit,

$$\nabla \ell_x^n(\theta) := \left( \frac{\partial}{\partial \theta_j} \ell_x^n(\theta) \right)_j = \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \log p_\theta(x_i) \right)_j.$$

Es gilt

$$\mathbb{E}_\theta[\nabla \ell_X^n(\theta)] = 0$$

wegen Proposition 5.21 sowie

$$n \text{COV}_\theta[\nabla \ell_X^n(\theta), \nabla \ell_X^n(\theta)] = I(\theta) = \left( -\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log p_\theta(X) \right] \right)_{jk}.$$

Damit folgt aus dem mehrdimensionalen Zentralen Grenzwertsatz bereits, dass für  $X \sim \mathbb{P}_\theta^\infty$

$$\sqrt{n} \nabla \ell_X^n(\theta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, I(\theta)). \quad (**)$$

Außerdem ist

$$\nabla \ell_x^n(d^n(x)) = 0,$$

da die Funktion  $\theta \mapsto \ell_x^n(\theta)$  bei  $d^n(x)$  ein Maximum besitzt und mit einer Taylor-Entwicklung von  $\nabla \ell_x^n(\eta)$  für  $\eta = d^n(x)$  um  $\theta$  gilt

$$0 = \nabla \ell_x^n(d^n(x)) = \nabla \ell_x^n(\theta) + B_n(d^n(x) - \theta)$$

für

$$B_n = \left( \frac{\partial^2}{\partial \eta_j \partial \eta_k} \ell_X^n(\eta) \Big|_{\eta = \eta_{n,j}^*} \right)_{jk}$$

für ein  $\eta_{n,j}^*$  zwischen  $\theta_j$  und  $d^n(x)_j$  für  $j = 1, \dots, m$ . Da  $d^n(X) \xrightarrow{n \rightarrow \infty}_p \theta$ , gilt auch  $\eta_{n,j}^* \xrightarrow{n \rightarrow \infty}_p \theta_j$ . Zusammen mit (\*\*) haben wir damit gezeigt, dass

$$\begin{aligned} & \sqrt{n} B_n(d^n(X) - \theta) \\ &= \sqrt{n} \left( \left( \frac{\partial^2}{\partial \eta_j \partial \eta_k} \ell_X^n(\eta) \Big|_{\eta = \eta_{n,j}^*} - \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X) \Big|_{\eta = \eta_{n,j}^*} \right] \right)_{jk} \right. \\ & \quad \left. + \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X) \Big|_{\eta = \eta_{n,j}^*} - \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log p_\theta(X) \right] \right)_{jk} \\ & \quad - I(\theta) \Big) (d^n(X) - \theta) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, I(\theta)). \end{aligned}$$

Nun gilt für  $X \sim \mathbb{P}_\theta^\infty$  wegen dem Gesetz der großen Zahlen und Annahme (\*)

$$\begin{aligned} & \frac{\partial^2}{\partial \eta_j \partial \eta_k} \ell_X^n(\eta) \Big|_{\eta = \eta_{n,j}^*} - \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X) \Big|_{\eta = \eta_{n,j}^*} \right] \xrightarrow{n \rightarrow \infty}_p 0, \\ & \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \eta_j \partial \eta_k} \log p_\eta(X) \Big|_{\eta = \eta_{n,j}^*} - \frac{\partial^2}{\partial \theta_j \partial \theta_k} \log p_\theta(X) \right] \xrightarrow{n \rightarrow \infty}_p 0. \end{aligned}$$



Damit folgt zunächst  $B_n \xrightarrow[n \rightarrow \infty]{p} -I(\theta)$ , woraus  $\sqrt{n}(d^n(X - \theta)) = O(1)$  folgt und damit (aus Slutskys Theorem)

$$-\sqrt{n}I(\theta)(d^n(X) - \theta) \xrightarrow[n \rightarrow \infty]{} \mathcal{N}(0, I(\theta)).$$

Da die Matrix-Multiplikation mit  $I(\theta)^{-1}$  eine stetige Abbildung ist, folgt damit das Ergebnis.  $\square$