

## 1

## Grundbegriffe der mathematischen Statistik

Elementare statistische Berechnungen werden schon seit Jahrtausenden durchgeführt. Das arithmetische Mittel aus einer Anzahl von Mess- oder Beobachtungswerten ist schon sehr lange bekannt.

Zuerst entstand die beschreibende Statistik mit dem Sammeln von Daten etwa bei Volkszählungen oder in Krankenregistern und deren Verdichtung in Form von Maßzahlen oder Grafiken. Die mathematische Statistik entwickelte sich ab Ende des 19. Jahrhunderts aufbauend auf der Wahrscheinlichkeitsrechnung. Anfang des 20. Jahrhunderts gehörten vor allem Karl Pearson und Sir Ronald Aymler Fisher zu ihren Pionieren. Das Buch von Fisher (1925) ist ein Meilenstein, in ihm werden die vom Autor mehrere Jahre zuvor entwickelten Grundlagen der Statistik wie die Maximum-Likelihood-Methode und die Varianzanalyse oder Begriffe wie Suffizienz und Effizienz Versuchsanstellern nahegebracht. Ein wichtiges Informationsmaß heißt noch heute Fisher-Information (siehe Abschn. 1.4).

Wir wollen auf die Details der historischen Entwicklung nicht eingehen und verweisen Interessierte auf Stigler (2000). Stattdessen beschreiben wir den heutigen Stand der Theorie. Wir wollen aber nicht vergessen, dass viele Anregungen aus Anwendungen kamen und bringen deshalb auch immer wieder Beispiele.

Die Wahrscheinlichkeitsrechnung ist zwar die Grundlage der mathematischen Statistik, aber viele praktische Probleme, in denen Aussagen über Zufallsvariablen gemacht werden sollen, sind mit der Wahrscheinlichkeitsrechnung allein nicht zu lösen. Das liegt daran, dass über die Verteilungsfunktion der Zufallsvariablen nicht alles bekannt ist und das Problem oft darin besteht, Aussagen über mindestens einen der Parameter einer Verteilungsfunktion zu machen oder dass sogar die Verteilungsfunktion gänzlich unbekannt ist. Die mathematische Statistik wird in vielen einführenden Texten als die Theorie der Auswertung von Versuchen oder Erhebungen betrachtet, d. h., man geht davon aus, dass bereits eine Zufallsstichprobe (nach Abschn. 1.1) vorliegt. Wie man auf optimalem Weg zu dieser Zufallsstichprobe gelangt, bleibt meist unberücksichtigt – dies wird gesondert in der statistischen Versuchsplanung abgehandelt. In den Anwendungen ist es klar, dass man erst den Versuch (die Erhebung) plant und dann, wenn der Versuch durchgeführt wurde, mit der Auswertung beginnt. In der Theorie ist es aber zweckmäßig, zunächst die optimale Auswertung zu ermitteln, um dann für diese den optimalen Versuchsplan zu bestimmen, z. B. den kleinsten Versuchsumfang für eine varianz-

optimale Schätzfunktion. Daher wird hier so verfahren, dass zunächst einmal die optimale Auswertung bestimmt wird und später für diese der Versuchsplan zu erarbeiten ist. Eine Ausnahme bilden dabei die sequentiellen Verfahren, bei denen Planung und Auswertung gemeinsam vorgenommen werden.

Wir müssen uns darüber im Klaren sein, dass es sich bei der Behandlung der mathematischen Statistik einerseits und bei ihrer Anwendung auf konkretes Datenmaterial andererseits um zwei völlig verschiedene Begriffssysteme handelt. In beiden treten oft die gleichen Termini auf, die es genau auseinanderzuhalten gilt. Wir sprechen davon, dass den Begriffen der empirischen Ebene (also denen der Realwelt) Modelle in der Theorie zugeordnet werden.

## 1.1

### Grundgesamtheit und Stichprobe

#### 1.1.1

##### Konkrete Stichproben und Grundgesamtheiten

In den empirischen Wissenschaften werden ein Merkmal oder auch mehrere Merkmale gleichzeitig (ein Merkmalsvektor) an bestimmten Objekten (oder Individuen) beobachtet. Aus den Beobachtungswerten sind Schlüsse auf die Gesamtheit der Merkmalswerte aller Objekte einer Gesamtheit zu ziehen. Ursache dafür ist, dass es sachliche oder ökonomische Gesichtspunkte gibt, die eine vollständige Erfassung der Merkmale aller Objekte nicht ermöglichen. Hierzu einige Beispiele:

- Die Kosten der Erfassung aller Merkmalswerte stehen in keinem Verhältnis zum Wert der Aussage (z. B. Messung der Körpergröße aller zurzeit lebenden Menschen über 18 Jahren).
- Die Erfassung der Merkmalswerte ist mit der Zerstörung der Objekte verbunden (nicht zerstörungsfreie Werkstoffprüfung wie Reißfestigkeit von Tauen oder Strümpfen).
- Die Gesamtheit der Objekte ist hypothetischer Natur, z. B. weil sie teilweise zum Untersuchungszeitpunkt nicht existieren (wie alle Produkte einer Maschine).

Die wenigen praktischen Fälle, in denen alle Objekte einer Gesamtheit beobachtet werden und auf keine umfassendere Gesamtheit geschlossen werden soll, können wir vernachlässigen, für sie benötigt man die mathematische Statistik nicht. Wir gehen also davon aus, dass aus einer Gesamtheit nur eine Teilmenge ausgewählt wird, um das Merkmal (den Merkmalsvektor) zu beobachten von dem auf die gesamte Population geschlossen werden soll. Einen solchen Teil nennen wir (konkrete) Stichprobe (der Objekte). Die Menge der an diesen Objekten gemessenen Merkmalswerte nennen wir (konkrete) Stichprobe der Merkmalswerte. Jedes Objekt der Population soll einen Merkmalswert besitzen (unabhängig davon, ob

wir ihn erfassen oder nicht). Die der Population entsprechende Gesamtheit der Merkmalswerte der Objekte dieser Population nennen wir Grundgesamtheit.

Eine Population und das zu erfassende Merkmal und damit auch die Grundgesamtheit müssen eindeutig definiert sein. Populationen sind vor allem räumlich und zeitlich abzugrenzen. Von einem beliebigen Objekt der Realwelt muss prinzipiell feststehen, ob es zur Population gehört oder nicht. Wir betrachten im Folgenden einige Beispiele:

Population		Grundgesamtheit	
A	Färsen einer bestimmten Rasse	$A_1$	Jahresmilchmenge dieser Färsen
	eines bestimmten Gebietes	$A_2$	180-Tage-Körpermasse dieser Färsen
	in einem bestimmten Jahr	$A_3$	Rückenhöhe dieser Färsen
B	Bewohner einer Stadt	$B_1$	Blutdruck dieser Bewohner um 6:00 Uhr
	an einem bestimmten Tag	$B_2$	Alter der Bewohner

Es ist einleuchtend, dass Schlüsse von der Stichprobe auf die Grundgesamtheit falsch sein können. Wenn man z. B. aus der Population  $B$  die Kinder einer Kindertagesstätte auswählt, ist möglicherweise der Blutdruck, aber ganz sicher das Alter nicht auf die Population verallgemeinerbar. Generell sprechen wir von Merkmalen, sofern diese aber einen bestimmten Einfluss auf die Versuchsergebnisse haben können, nennen wir sie auch Faktoren, die (meist wenigen) Merkmalswerte heißen dann Faktorstufen, die Kombination von Faktorstufen mehrerer Faktoren heißen Faktorstufenkombinationen.

Hinsichtlich aller Faktoren, die das Merkmal in einer Grundgesamtheit beeinflussen können, sollte die Stichprobe „repräsentativ“ sein. Das heißt, in der Stichprobe der Objekte sollte sich die Zusammensetzung der Population widerspiegeln. Das ist aber bei kleinen Stichproben und vielen Faktorstufenkombinationen gar nicht möglich. In Population  $B$  gibt es hinsichtlich der Faktoren Alter und Geschlecht schon etwa 200 Faktorstufenkombinationen, die sich unmöglich in einer Stichprobe von 100 Einwohnern widerspiegeln können. Wir empfehlen daher, den Begriff „repräsentative Stichprobe“ nicht zu verwenden, da er nicht sauber definiert werden kann.

Stichproben sollen nicht danach beurteilt werden, welche Elemente sie enthalten, sondern danach, wie sie erhalten (gezogen) wurden. Die Art und Weise, wie eine Stichprobe erhoben wird, heißt Stichprobenverfahren. Es kann entweder auf die Objekte als Merkmalsträger oder auf die Grundgesamtheit der Merkmalswerte (z. B. in einer Datenbank) angewendet werden. Im letzteren Fall entsteht die Stichprobe der Merkmalswerte unmittelbar. Im ersteren Fall muss das Merkmal an den ausgewählten Objekten noch erfasst werden. Beide Vorgehensweisen (nicht unbedingt die entstehenden Stichproben) sind dann identisch, wenn für jedes ausgewählte Objekt der Merkmalswert erfasst wird. Davon gehen wir in diesem Kapitel aus. In zensierten Stichproben ist das nicht der Fall. Eine Stichprobe heißt zensiert, wenn der Merkmalswert nicht an allen Versuchseinheiten erfasst

werden konnte. Bricht man z. B. eine Lebensdauerermittlung (z. B. von elektronischen Bauteilen) nach einer bestimmten Zeit ab, liegen Messwerte für Objekte mit längerer Lebensdauer (als die Beobachtungszeit) nicht vor.

Im Folgenden wird nicht zwischen Stichproben der Objekte und der Merkmalswerte unterschieden, die Definitionen gelten für beide.

### Definition 1.1

Ein Stichprobenverfahren ist eine Vorschrift für die Auswahl einer endlichen Teilmenge, genannt Stichprobe, aus einer wohldefinierten endlichen Population (Grundgesamtheit), es heißt zufällig, wenn jedes Element der Grundgesamtheit mit der gleichen Wahrscheinlichkeit  $p$  in die Stichprobe gelangen kann. Eine (konkrete) Stichprobe ist das Ergebnis der Anwendung eines Stichprobenverfahrens. Stichproben, die das Ergebnis eines zufälligen Stichprobenverfahrens sind, heißen (konkrete) zufällige Stichproben oder (konkrete) Zufallsstichproben.

In der Stichprobentheorie (siehe z. B. Cochran und Boing 1972; Kauermann und Küchenhoff 2011 oder Quatember 2014) wird eine Vielzahl von zufälligen Stichprobenverfahren zur Verfügung gestellt. Wir verwenden ab jetzt die Begriffe Population und Grundgesamtheit synonym.

#### 1.1.2

### Stichprobenverfahren

Bei Zufallsauswahlverfahren unterscheiden wir u. a.:

- die einfache oder reine Zufallsauswahl, bei der jedes Element der Grundgesamtheit die gleiche Wahrscheinlichkeit hat, in die Stichprobe zu gelangen.
- die geschichtete Auswahl, bei der innerhalb zuvor festgelegter (disjunkter) Klassen eine zufällige Auswahl vorgenommen wird, sie ist nur dann insgesamt zufällig, wenn die Auswahlwahrscheinlichkeiten innerhalb der Klassen proportional zum Umfang der Klassen gewählt werden.
- die Klumpenauswahl, hier wird eine Grundgesamtheit in Gruppen (Klumpen) eingeteilt. Die Auswahl der Untersuchungsobjekte erfolgt nicht unter den Elementen der Grundgesamtheit, sondern unter den (disjunkten) Klumpen. In den ausgewählten Klumpen werden dann alle Elemente erfasst. Sie findet häufig in Form der Flächenstichproben Anwendung. Sie ist nur dann zufällig im Sinne der Definition 1.1, wenn die Klumpen gleich viele Elemente enthalten.
- die mehrstufige Auswahl, sie ist dadurch gekennzeichnet, dass mindestens zwei Auswahlstufen bestehen. Die Grundgesamtheit wird z. B. in zweistufiger Auswahl in Primäreinheiten in Form disjunkter Teilmengen zerlegt. Aus der Menge der Primäreinheiten erfolgt zunächst eine Zufallsauswahl. Aus jeder ausgewählten Primäreinheit erfolgt eine Zufallsauswahl von Untersuchungseinheiten (Sekundäreinheiten). Eine mehrstufige Auswahl ist dann vorteilhaft, wenn die Grundgesamtheit hierarchisch gegliedert ist (Land, Provinzen, Städ-

te in der Provinz). Sie ist nur dann zufällig im Sinne der Definition 1.1, wenn die Primäreinheiten gleich viele Sekundäreinheiten enthalten.

- die stets sequentielle Auswahl; hier liegt der Stichprobenumfang nicht vor Beginn des Auswahlprozesses fest, es wird zunächst eine kleine Stichprobe gezogen und analysiert. Es erfolgt dann eine Entscheidung, ob die vorliegende Information hinreichend ist, z. B. um eine Hypothese abzulehnen oder anzunehmen (siehe Kapitel 3), oder ob mehr Information durch Ziehung einer weiteren Einheit beschafft werden soll.

Sowohl ein zufälliges Stichprobenverfahren als auch eine willkürliche Auswahl aufs Geratewohl können zur gleichen konkreten Stichprobe führen. Ob einer konkreten Stichprobe eine zufällige oder eine willkürliche Auswahl zugrunde liegt, kann nicht anhand dieser Stichprobe beurteilt werden, sondern eben nur anhand des verwendeten Auswahlverfahrens.

Bei der reinen Zufallsauswahl wird Definition 1.1 direkt angewendet, jedes Element einer Grundgesamtheit vom Umfang  $N$  wird mit der gleichen Wahrscheinlichkeit  $p$  der Grundgesamtheit entnommen. Wir nennen die Anzahl der Elemente in einer Stichprobe den Stichprobenumfang und bezeichnen diesen in der Regel mit  $n$ .

Der praktisch wichtige Fall einer reinen Zufallsauswahl ist der, dass entnommene Elemente nicht in die Grundgesamtheit zurückgelegt werden, wie das etwa bei der Ziehung der Lottozahlen der Fall ist. Hier werden in Deutschland  $n = 6$  Zahlen aus  $N = 49$  gegebenen Zahlen gezogen. Bei einer uneingeschränkten Zufallsstichprobe vom Umfang  $n$  haben alle möglichen  $\binom{N}{n}$  Teilmengen die gleiche Wahrscheinlichkeit  $\frac{1}{\binom{N}{n}}$  dafür, in die Stichprobe zu gelangen.

Ob eine Stichprobe eine Zufallsstichprobe ist oder nicht, kann man ihr – wie gesagt – nicht ansehen. Man muss vielmehr das Verfahren betrachten, mit dem sie gezogen wurde. Allerdings wird man sofort misstrauisch, wenn extreme Stichproben auftreten. Wird aus einer Grundgesamtheit mit 10 000 Losen der Hauptgewinn während des Kaufs eines speziellen Loses gezogen, so ist das schon ungewöhnlich, kann aber, wie man im Volksmund sagt, schon mit rechten Dingen zugegangen sein, in unserer Terminologie also das Ergebnis eines Zufallsstichprobenverfahrens sein. Zieht dieselbe Person an drei aufeinanderfolgenden Verlosungen den Hauptgewinn und stellt sich dann noch heraus, dass es sich um den Bruder des Losverkäufers handelt, stellen sich berechtigte Zweifel ein. Wir weigern uns, Ereignisse mit solch geringer Wahrscheinlichkeit zu akzeptieren und vermuten, dass das zugrunde gelegte Modell falsch ist. In diesem Fall nehmen wir an, dass kein Zufallsstichprobenverfahren zugrunde gelegt wurde und Betrug im Spiel ist. Trotzdem besteht eine ganz geringe Wahrscheinlichkeit für dieses Ereignis als Ergebnis eines Zufallsstichprobenverfahrens, nämlich  $1/1\,000\,000\,000\,000$ .

Nebenbei bemerkt bildet diese Art, Modelle (Sachverhalte) zu verwerfen, unter denen ein beobachtetes Ereignis eine sehr kleine Wahrscheinlichkeit hat und statt dessen solche Modelle zu akzeptieren, bei denen die Wahrscheinlichkeit dieses Ereignisses größer ist, die Basis für die statistischen Tests in Kapitel 3.

Bei einer Zufallsstichprobe mit Zurücklegen wird auch ein reines Stichprobenverfahren verwendet, also jedes Element hat die gleiche Wahrscheinlichkeit mit Zurücklegen gezogen zu werden. Es wird jedes gezogene und beobachtete Element in die Grundgesamtheit zurückgelegt, bevor das nächste Element gezogen wird. Das geht nur bei zerstörungsfreier Beobachtung, d. h. in solchen Fällen, bei denen sich die Stichprobeneinheit durch die Beobachtung nicht verändert (Beispiele, bei denen ein Zurücklegen nicht möglich ist, sind Zerreißproben, Untersuchungen an geschlachteten Tieren, Fällen von Bäumen, Abernten u. a.). Dieses Verfahren heißt einfaches Zufallsstichprobenverfahren mit Zurücklegen. Bei Zufallsstichprobenverfahren ohne Zurücklegen erhält man  $n < N$  verschiedene Elemente, beim Zufallsstichprobenverfahren mit Zurücklegen kann dasselbe Element mehrfach in der Stichprobe auftauchen und es kann auch  $n > N$  sein.

Eine mitunter praktisch einfacher zu realisierende Methode ist die systematische Auswahl mit Zufallsstart. Sie ist anwendbar, wenn die Elemente der endlichen Auswahlgrundlage von 1 bis  $N$  durchnummeriert sind und die Folge nicht mit dem Merkmal zusammenhängt. Wenn  $N/n$  ganzzahlig ist, wählt man zufällig eine Zahl  $i$  zwischen 1 und  $N/n$  aus und bildet die Stichprobe aus den Elementen  $i, N/n + i, 2N/n + i, \dots, (n-1)N/n + i$ . Näheres hierzu und was zu tun ist, wenn  $N/n$  nicht ganzzahlig ist, findet man bei (Rasch *et al.*, 2008), Verfahren 1/31/1210.

Die oben erwähnte geschichtete Auswahl bietet sich dann an, wenn die Grundgesamtheit vom Umfang  $N$  auf inhaltlich relevante Weise in  $s$  Teilgesamtheiten vom Umfang  $N_1, N_2, \dots, N_s$  zerfällt. Insbesondere kann die Grundgesamtheit gelegentlich nach den Stufen eines vermuteten Störfaktors in solche Teilgesamtheiten unterteilt werden. Man bezeichnet diese Teilgesamtheiten als Schichten. Will man aus dieser Grundgesamtheit Stichproben vom Umfang  $n$  erheben, so muss man bei einem uneingeschränkten Zufallsstichprobenverfahren befürchten, dass nicht alle Schichten überhaupt bzw. zumindest nicht in angemessener Weise berücksichtigt werden. Dann ist es besser, ein geschichtetes Zufallsstichprobenverfahren durchzuführen. Man erhebt dabei jeweils Teilstichproben vom Umfang  $n_i$  ( $i = 1, 2, \dots, s$ ) aus der  $i$ -ten Schicht. Die Teilstichproben werden aus der jeweiligen Schicht nach einem reinen Zufallsstichprobenverfahren gezogen. Dies entspricht, wenn  $n_i/n$  proportional zu  $N_i/N$  gewählt wird, auch insgesamt einem Zufallsstichprobenverfahren.

Während beim geschichteten Zufallsstichprobenverfahren aus jeder Teilmenge Elemente erhoben werden, werden bei der mehrstufigen Auswahl wie oben beschrieben auf jeder Stufe zufällig Teilmengen oder Elemente entnommen. Im zweistufigen Fall bestehe die Grundgesamtheit aus  $k$  disjunkten Teilmengen vom Umfang  $N_0$ , den Primäreinheiten. Es wird nun vorausgesetzt, dass sich die Merkmalswerte zwischen den Primäreinheiten nur zufällig unterscheiden, sodass nicht aus allen Primäreinheiten Elemente entnommen werden müssen. Ist der gewünschte Stichprobenumfang  $n = rn_0$  mit  $r < k$ , so wählt man zunächst nach einem reinen Zufallsstichprobenverfahren  $r$  der  $k$  Primäreinheiten aus. Aus jeder der  $r$  in der ersten Stufe erhobenen Primäreinheiten wählt man in der zweiten Stufe nach einem reinen Zufallsstichprobenverfahren je  $n_0$  Objekte (Sekundär-

Tab. 1.1 Anzahl  $K$  möglicher Stichproben für verschiedene Stichprobenverfahren.

Stichprobenverfahren	Anzahl $K$ möglicher Stichproben
Reine Zufallsauswahl	$K = \binom{1000}{100} > 10^{140}$
Systematische Auswahl mit Zufallsstart $k = 10$	$K = 10$
Geschichtetes Stichprobenverfahren, $k = 20, N_i = 50, i = 1, \dots, 20$	$K = \left[ \binom{50}{5} \right]^{20} = 2,118\,76 \cdot 10^{120}$
Geschichtetes Stichprobenverfahren, $k = 10, N_i = 100, i = 1, \dots, 10$	$K = \left[ \binom{100}{10} \right]^{10} = 1,731\,030\,9 \cdot 10^{130}$
Geschichtetes Stichprobenverfahren, $k = 5, N_i = 200, i = 1, \dots, 5$	$K = \left[ \binom{200}{20} \right]^5 = 1,613\,587\,8 \cdot 10^{135}$
Geschichtetes Stichprobenverfahren, $k = 2, N_1 = 400, N_2 = 600$	$K = \binom{400}{40} \cdot \binom{600}{60} = 5,466\,241\,4 \cdot 10^{138}$
Zweistufiges Verfahren $k = 20, N_0 = 50, r = 4$	$K = \binom{20}{4} \cdot \binom{50}{25} = 6,124\,593\,9 \cdot 10^{17}$
Zweistufiges Verfahren $k = 20, N_0 = 50, r = 5$	$K = \binom{20}{5} \cdot \binom{50}{20} = 7,306\,913\,1 \cdot 10^{17}$
Zweistufiges Verfahren $k = 10, N_0 = 100, r = 2$	$K = \binom{10}{2} \cdot \binom{100}{50} = 4,540\,110\,5 \cdot 10^{30}$
Zweistufiges Verfahren $k = 10, N_0 = 100, r = 4$	$K = \binom{10}{4} \cdot \binom{100}{25} = 5,092\,904\,7 \cdot 10^{25}$
Zweistufiges Verfahren $k = 5, N_0 = 200, r = 2$	$K = \binom{5}{2} \cdot \binom{200}{50} = 4,538\,583\,8 \cdot 10^{48}$
Zweistufiges Verfahren $k = 2, N_0 = 500, r = 1$	$K = \binom{2}{1} \cdot \binom{500}{100} > 10^{100}$

einheiten) aus. Die Anzahl der möglichen Stichproben beträgt  $\binom{k}{r} \cdot \binom{N_0}{n_0}$  und entsprechend Definition 1.1 gelangt jedes Element der Grundgesamtheit mit der gleichen Wahrscheinlichkeit  $p = \frac{r}{k} \cdot \frac{n_0}{N_0}$  in die Stichprobe.

### Beispiel 1.1

Aus einer Grundgesamtheit mit  $N = 1000$  Objekten soll eine Zufallsstichprobe ohne Zurücklegen vom Umfang  $n = 100$  gezogen werden. Tabelle 1.1 gibt für verschiedene Verfahren die Anzahl der Stichproben an, die Wahrscheinlichkeit der Auswahl ist für jedes Objekt  $p = 0,1$ .

## 1.2

### Mathematische Modelle für Grundgesamtheit und Stichprobe

In der mathematischen Statistik werden Begriffe definiert, die als Modelle (Verallgemeinerungen) für die in der Empirie gebräuchlichen Begriffe verwendet wer-

den. Der Grundgesamtheit, die einer Häufigkeitsverteilung der Merkmalswerte entspricht, wird als Modell die Wahrscheinlichkeitsverteilung gegenübergestellt. Die durch ein Zufallsstichprobenverfahren entstandene konkrete Stichprobe wird durch die realisierte (theoretische) Zufallsstichprobe modelliert. Diese Modellvorstellungen sind dann adäquat, wenn der Umfang  $N$  der Grundgesamtheit sehr groß im Vergleich zum Umfang  $n$  der Stichprobe ist.

### Definition 1.2

Eine  $n$ -dimensionale Zufallsvariable

$$Y = (y_1, y_2, \dots, y_n)^T, \quad n \geq 1$$

mit den Komponenten  $y_i$  heißt Zufallsstichprobe, wenn

- die  $y_i$  die gleiche durch die Verteilungsfunktion  $F(y_i, \theta) = F(y, \theta)$  charakterisierte Verteilung mit dem Parameter(vektor)  $\theta \in \Omega \subseteq R^p$  haben und
- die  $y_i$  voneinander stochastisch unabhängig sind, sodass für die Verteilungsfunktion  $F(Y, \theta)$  von  $Y$

$$F(Y, \theta) = \prod_{i=1}^n F(y_i, \theta), \quad \theta \in \Omega \subseteq R^p$$

gilt.

Die Werte  $Y = (y_1, y_2, \dots, y_n)^T$  einer Zufallsstichprobe  $Y$  heißen Realisationen. Die Gesamtheit  $\{Y\}$  aller möglichen Realisationen von  $Y$  heißt Stichprobenraum.

In diesem Buch werden Zufallsvariablen fett gedruckt, und der Stichprobenraum  $\{Y\}$  liegt stets im  $n$ -dimensionalen euklidischen Raum, d. h.  $\{Y\} \subset R^n$ .

Die Funktion

$$L(Y, \theta) = \begin{cases} f(Y, \theta) = \frac{\partial F(Y, \theta)}{\partial Y}, & \text{für kontinuierliche } y \\ p(Y, \theta), & \text{für diskrete } y \end{cases}$$

mit der Wahrscheinlichkeitsfunktion  $p(Y, \theta)$  bzw. der Dichtefunktion  $f(Y, \theta)$  bei gegebenem  $Y$  als Funktion von  $\theta$  heißt Likelihood-Funktion.

Das Wort Zufallsstichprobe kann nun folgendes bedeuten:

- Zufallsstichprobe als Zufallsvariable  $Y$  nach Definition 1.2,
- (konkrete) Zufallsstichprobe als Teilmenge einer Population (Grundgesamtheit), die nach einem Zufallsstichprobenverfahren ausgewählt wurde.

Die Realisationen  $Y$  einer Zufallsstichprobe  $Y$  werden wir dagegen stets realisierte Zufallsstichprobe nennen.

Eine Zufallsstichprobe  $Y$  ist das mathematische Modell des reinen Zufallsstichprobenverfahrens, konkrete Zufallsstichprobe und realisierte Zufallsstichprobe entsprechen einander auch in der Symbolik.



Wir beschreiben in diesem Buch das „klassische“ Vorgehen, nach dem  $Y$  mit der Verteilungsfunktion  $F(Y, \theta)$  mit dem festen (nicht zufälligen) Parameter  $\theta \in \Omega \subseteq R^p$  verteilt ist. Daneben gibt es das Bayessche Vorgehen, bei dem man ein zufälliges  $\theta$  annimmt, das mit einer a-priori-Verteilung mit einem als bekannt vorausgesetztem Parameter  $\varphi$  verteilt ist. Beim empirischen Bayesschen Vorgehen wird die a-priori-Verteilung aus bereits ermittelten Daten geschätzt.

### 1.3

#### Suffizienz und Vollständigkeit

Eine Zufallsgröße enthält gewisse Informationen über die Verteilung und deren Parameter. Vor allem für große  $n$  (etwa  $n > 100$ ) möchte man die Elemente der Zufallsstichprobe so verdichten, dass möglichst wenige neue Zufallsvariablen möglichst viel von dieser Information enthalten. Diese unklar formulierte Wunschvorstellung soll jetzt schrittweise bis zum Konzept der minimal suffizienten Maßzahl präzisiert werden. Zunächst wiederholen wir hier die Definition einer Exponentialfamilie.

Die Verteilung einer Zufallsvariablen  $y$  mit dem Parametervektor  $\theta = (\theta_1, \theta_2, \dots, \theta_p)^T$  gehört zu einer  $k$ -parametrischen Exponentialfamilie, wenn ihre Likelihood-Funktion in der Form

$$f(y, \theta) = h(y)e^{\sum_{i=1}^k \eta_i(\theta) \cdot T_i(y) - B(\theta)}$$

geschrieben werden kann, wobei folgendes gilt:

- $\eta_i$  und  $B$  sind reelle Funktionen von  $\theta$  und  $B$  hängt nicht von  $y$  ab.
- Die Funktion  $h(y)$  ist nichtnegativ und hängt nicht von  $\theta$  ab.

Die Exponentialfamilie ist in kanonischer Form mit den sogenannten natürlichen Parametern  $\eta_i$ , falls ihre Elemente als

$$f(y, \eta) = h(y)e^{\sum_{i=1}^k \eta_i \cdot T_i(y) - A(\eta)} \quad \text{mit} \quad \eta = (\eta_1, \dots, \eta_k)^T$$

geschrieben werden können.

Wir gehen von einer Verteilungsfamilie  $(P_\theta, \theta \in \Omega)$  von Zufallsvariablen  $y$  mit der Verteilungsfunktion  $F(y, \theta)$ ,  $\theta \in \Omega$  aus. Die Realisationen  $Y = (y_1, \dots, y_n)^T$  der Zufallsstichprobe

$$Y = (y_1, y_2, \dots, y_n)^T$$

mit wie  $y$  verteilten Komponenten liegen im Stichprobenraum  $\{Y\}$ . Nach Definition 1.2 ist mit  $F(y, \theta)$  auch die Verteilungsfunktion  $F(Y, \theta)$  einer Zufallsstichprobe  $Y$  eindeutig festgelegt.

**Definition 1.3**

Eine messbare Abbildung  $M = M(Y) = [M_1(Y), \dots, M_r(Y)]^T$ ,  $r \leq n$  von  $\{Y\}$  auf einen Raum  $\{M\}$ , die nicht von  $\theta \in \Omega$  abhängt, heißt (statistische) Maßzahl oder auch Statistik.

**Definition 1.4**

Eine Maßzahl  $M$  heißt *suffizient* oder *erschöpfend* bezüglich einer Verteilungsfamilie  $(P_\theta, \theta \in \Omega)$  bzw. bezüglich  $\theta \in \Omega$ , falls die bedingte Verteilung einer Zufallsstichprobe  $Y$  bei gegebenem  $M = M(Y) = M(Y)$  von  $\theta$  unabhängig ist.

**Beispiel 1.2**

Die Komponenten einer Zufallsstichprobe  $Y$  mögen einer Zweipunktverteilung mit den Werten 1 und 0 folgen. Dabei sei  $P(y_i = 1) = p$  und  $P(y_i = 0) = 1 - p$  mit  $0 < p < 1$ . Dann ist  $M = M(Y) = \sum_{i=1}^n y_i$  suffizient bezüglich  $\theta \in (0, 1) = \Omega$ . Um das zu zeigen, müssen wir nachweisen, dass  $P(Y = Y | \sum_{i=1}^n y_i = M)$  von  $p$  unabhängig ist. Nun ist  $P(Y = Y | M) = \frac{P(Y=Y, M=M)}{P(M=M)}$ ,  $M = 0, 1, \dots, n$ . Aus der Wahrscheinlichkeitsrechnung wissen wir, dass  $M = M(Y) = \sum_{i=1}^n y_i$  binomialverteilt ist mit den Parametern  $n$  und  $p$ , also gilt:

$$P(M = M) = \binom{n}{M} p^M (1-p)^{n-M}, \quad M = 0, 1, \dots, n$$

Ferner ist mit  $y_i = 0$  oder  $y_i = 1$  und  $A(M) = \{Y | M(Y) = M\}$

$$\begin{aligned} P[Y = Y, M(Y) = M] &= P(\mathbf{y}_1 = y_1, \dots, \mathbf{y}_n = y_n) I_{A(M)}(Y) \\ &= \prod_{i=1}^n \left( \binom{1}{y_i} p^{y_i} (1-p)^{1-y_i} \right) I_{A(M)}(Y) \\ &= p^{\sum_{i=1}^n y_i} (1-p)^{n - \sum_{i=1}^n y_i} I_{A(M)}(Y) \\ &= p^M (1-p)^{n-M} \end{aligned}$$

Daher ist  $P(Y = Y | M) = \frac{1}{\binom{n}{M}}$  und das ist unabhängig von  $p$ .

Auf diese Weise ist der Nachweis der Suffizienz recht mühsam, er gelingt aber auch für kontinuierliche Verteilungen, wie das nächste Beispiel zeigt.

**Beispiel 1.3**

Die Komponenten  $y_i$  einer Zufallsstichprobe  $Y$  vom Umfang  $n$  seien nach  $N(\mu, 1)$  mit Erwartungswert  $\mu$  und Varianz  $\sigma^2 = 1$  verteilt. Dann ist  $M = \sum y_i$  suffizient bezüglich  $\mu \in \mathbb{R}^1 = \Omega$ . Um das zu zeigen, vermerken wir zunächst, dass  $Y$  nach  $N(\mu e_n, E_n)$  verteilt ist. Nun führen wir die eineindeutige Transformation

$$Z = AY = (z_1 = \sum y_i, z_2 = y_1, \dots, z_n = y_1) \quad \text{mit} \quad A = \begin{pmatrix} 1 & \mathbf{e}_{n-1}^T \\ -e_{n-1} & E_{n-1} \end{pmatrix}$$

durch, es gilt  $|A| = n$ . Wir schreiben  $Z = (z_1, Z_2) = (\sum y_i, y_2 - y_1, \dots, y_n - y_1)$  und sehen, dass

$$\text{cov}(Z_2, z_1) = \text{cov}(Z_2, M) = \text{cov}((-e_{n-1}, E_{n-1})Y, e_n^T Y) = 0_{n-1}$$

gilt.

Wegen der Normalverteilungsannahme sind damit  $M$  und  $Z_2$  stochastisch unabhängig. Damit sind  $Z_2$ , aber auch  $Z_2|M$  und auch  $Z|M$  von  $\mu$  unabhängig. Wegen der Eineindeutigkeit der Abbildung  $Z = AY$  ist auch  $Y|M$  von  $\mu$  unabhängig und damit ist  $M = \sum y_i$  suffizient bezüglich  $\mu \in R^1$ . Mit  $M = \sum y_i$  und einer reellen Zahl  $c \neq 0$  ist stets auch  $cM$  also z. B.  $\frac{1}{n} \sum y_i = \bar{y}$  suffizient.

Die Suffizienz spielt nun aber in der mathematischen Statistik eine so große Rolle, dass wir einfachere Methoden zum Nachweis der Suffizienz und vor allem zum Auffinden suffizienter Maßzahlen benötigen. Der nachfolgende Satz hilft uns da weiter.

### Satz 1.1 Zerlegungssatz

Gegeben sei eine Verteilungsfamilie  $(P_\theta, \theta \in \Omega)$  einer Zufallsstichprobe  $Y$ , die von einem endlichen Maß  $\nu$  dominiert wird. Die Maßzahl  $M(Y)$  ist genau dann bezüglich  $\theta$  suffizient, wenn die Radon-Nikodym-Dichte  $f_\theta$  von  $P_\theta$  bezüglich  $\nu$  in der Form

$$f_\theta(Y) = g_\theta[M(Y)]h(Y) \quad (1.1)$$

$\nu$ -fast überall geschrieben werden kann, wobei gilt: die  $\nu$ -integrierbare Funktion  $g_\theta$  ist nichtnegativ und messbar,  $h$  ist nichtnegativ und  $h(Y) = 0$  nur für eine  $P_\theta$ -Nullmenge.

Der allgemeine Beweis stammt von Halmos und Savage (1949), man findet ihn auch z. B. bei Bahadur (1955) oder Lehmann (1959).

Wir beschäftigen uns in diesem Buch nur mit diskreten und kontinuierlichen Wahrscheinlichkeitsverteilungen, die die Voraussetzungen dieses Satzes erfüllen. Den Beweis des Satzes für solche Verteilungen gibt Rasch (1995). Wir verzichten hier auf dessen Wiederholung.

Für diskrete Verteilungen bedeutet dieser Satz, dass die Wahrscheinlichkeitsfunktion die Form

$$p(Y, \theta) = g[M(Y), \theta]h(Y) \quad (1.2)$$

hat. Für kontinuierliche Verteilungen hat die Dichtefunktion die Form

$$f(Y, \theta) = g[M(Y), \theta]h(Y) \quad (1.3)$$

**Korollar 1.1**

Ist die Verteilungsfamilie  $(P^*(\theta), \theta \in \Omega)$  der Zufallsvariablen  $\mathbf{y}$  eine  $k$ -parametrische Exponentialfamilie mit natürlichem Parameter  $\eta$  und der Likelihood-Funktion

$$L^*(\mathbf{y}, \eta) = h^*(\mathbf{y})e^{\sum_{j=1}^k M_j^*(\mathbf{y}) - A(\eta)} \quad (1.4)$$

so ist mit der Zufallsstichprobe  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)^T$

$$M(\mathbf{Y}) = \left( \sum_{i=1}^n M_1^*(\mathbf{y}_i), \dots, \sum_{i=1}^n M_k^*(\mathbf{y}_i) \right)^T \quad (1.5)$$

suffizient bezüglich  $\theta$ .

*Beweis:* Es gilt

$$L(\mathbf{y}, \eta) = \prod_{i=1}^n h^*(\mathbf{y}_i) e^{\sum_{j=1}^k \eta_j \sum_{i=1}^n M_j^*(\mathbf{y}_i) - nA(\eta)} \quad (1.6)$$

und das hat die Form (1.2) bzw. (1.3) mit  $h(\mathbf{Y}) = \prod_{i=1}^n h^*(\mathbf{y}_i)$  und  $\theta = \eta$ .

**Definition 1.5**

Zwei Likelihood-Funktionen  $L_1(Y_1, \theta)$  und  $L_2(Y_2, \theta)$  heißen äquivalent,  $L_1 \sim L_2$ , wenn

$$L_1(Y_1, \theta) = a(Y_1, Y_2)L_2(Y_2, \theta) \quad (1.7)$$

mit einer von  $\theta$  unabhängigen Funktion  $a(Y_1, Y_2)$  ist.

Dann folgt aus Satz 1.1

**Korollar 1.2**

$M(\mathbf{Y})$  ist genau dann suffizient bezüglich  $\theta$ , wenn die Likelihood-Funktion  $L_M(M, \theta)$  von  $\mathbf{M} = M(\mathbf{Y})$  äquivalent zur Likelihood-Funktion einer Zufallsstichprobe  $\mathbf{Y}$  ist.

*Beweis:* Ist  $M(\mathbf{Y})$  suffizient, so hat mit  $L(\mathbf{Y}, \eta)$  wegen

$$L_M(M, \theta) = a(\mathbf{Y})L(\mathbf{Y}, \theta), \quad a(\mathbf{Y}) > 0 \quad (1.8)$$

auch  $L_M(M, \theta)$  die Form (1.1). Gilt andererseits (1.8), so folgt, dass die bedingte Verteilung einer Zufallsstichprobe  $\mathbf{Y}$  bei gegebenem  $M(\mathbf{Y}) = M$  von  $\theta$  unabhängig ist.

**Beispiel 1.4**

Die Komponenten  $y_i$  einer Zufallsstichprobe  $Y = (y_1, y_2, \dots, y_n)^T$  seien nach  $N(\mu, 1)$  verteilt. Es gilt:

$$L(Y, \mu) = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2}(Y - \mu e_n)^T (Y - \mu e_n)} = \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{1}{2} \sum_{i=1}^n (y_i - \bar{y})^2} e^{-\frac{n}{2} (\bar{y} - \mu)^2} \quad (1.9)$$

Da  $M(Y) = \bar{y}$  nach  $N(\mu, \frac{1}{n})$  verteilt ist, ist

$$L_M(\bar{y}, \mu) = \frac{\sqrt{n}}{\sqrt{2\pi}} e^{-\frac{n}{2} (\bar{y} - \mu)^2} \quad (1.10)$$

und damit gilt  $L_M(\bar{y}, \mu) \sim L(Y, \mu)$  und  $\bar{y}$  ist suffizient bezüglich  $\mu$ .

Allgemein folgt unmittelbar aus Definition 1.4

**Korollar 1.3**

Ist  $c > 0$  eine von  $\theta$  unabhängig gewählte reelle Zahl und  $M(Y)$  suffizient bezüglich  $\theta$ , so ist auch  $cM(Y)$  suffizient bezüglich  $\theta$ .

So ist also z. B. mit  $M = \sum y_i$  und  $c = \frac{1}{n}$  auch  $\frac{1}{n} \sum y_i = \bar{y}$  suffizient.

Man kann nun die Frage stellen, ob es unter den suffizienten Maßzahlen bezüglich einer Verteilungsfamilie  $P^*(\theta)$ ,  $\theta \in \Omega$  solche Maßzahlen gibt, die in einem noch zu definierenden Sinne minimal sind, also möglichst wenige Komponenten enthalten. Wie das folgende Beispiel zeigt, ist diese Frage nicht abwegig.

**Beispiel 1.5**

Es sei  $P^*(\theta)$ ,  $\theta \in \Omega$  die Familie der  $N(\mu, \sigma^2)$ -Normalverteilungen ( $\sigma > 0$ ). Wir betrachten die Maßzahlen einer Zufallsstichprobe  $Y$  vom Umfang  $n$ :

$$\begin{aligned} M_1(Y) &= Y \\ M_2(Y) &= (y_1^2, \dots, y_n^2)^T \\ M_3(Y) &= \left( \sum_{i=1}^r y_i^2, \sum_{i=r+1}^n y_i^2 \right)^T, \quad r = 1, \dots, n-1 \\ M_4(Y) &= \left( \sum_{i=1}^n y_i^2 \right) \end{aligned}$$

die alle bezüglich  $\sigma^2$  suffizient sind. Das zeigt man sehr einfach mithilfe von Korollar 1.1 zum Zerlegungssatz. Die Likelihood-Funktion von  $M_1(Y)$  und  $Y$  sind identisch (und damit äquivalent). Da mit den  $y_i$  auch die  $y_i^2$  unabhängig sind und

die  $\frac{y_i^2}{\sigma^2} = \chi_i^2$  nach  $CQ(1)$   $\chi^2$ -verteilt sind, folgt nach der Transformation  $y_i^2 = \sigma^2 \chi_i^2$

$$L_M(M_2(\mathbf{Y}), \sigma^2) \sim L(\mathbf{Y}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2} \quad (1.11)$$

Analog verfährt man mit  $M_3(\mathbf{Y})$  und  $M_4(\mathbf{Y})$ .

Sicher stellt  $M_4(\mathbf{Y})$  die weitestgehende Zusammenfassung der Komponenten einer Zufallsstichprobe  $\mathbf{Y}$  dar und ist gegenüber den anderen Maßzahlen vorzuziehen.

### Definition 1.6

Eine bezüglich  $\theta$  suffiziente Maßzahl  $M^*(\mathbf{Y})$  heißt minimal suffizient bezüglich  $\theta$ , wenn sie sich als eine Funktion jeder anderen suffizienten Maßzahl  $M(\mathbf{Y})$  darstellen lässt.

Betrachten wir Beispiel 1.5, so ist

$$M_4(\mathbf{Y}) = M_1^T(\mathbf{Y})M_1(\mathbf{Y}) = e_n^T M_2 = \begin{pmatrix} 1 & 1 \end{pmatrix} M_3, \quad r = 1, \dots, n-1$$

Damit kann  $M_4(\mathbf{Y})$  als Funktion aller suffizienten Maßzahlen des Beispiels geschrieben werden. Für  $M_1(\mathbf{Y})$ ,  $M_2(\mathbf{Y})$  und  $M_3(\mathbf{Y})$  gilt das nicht, sie sind keine Funktionen von  $M_4(\mathbf{Y})$ .  $M_4(\mathbf{Y})$  ist die einzige Maßzahl des Beispiels 1.5, die minimal suffizient bezüglich  $\sigma^2$  sein könnte. Wir werden sehen, dass sie tatsächlich diese Eigenschaft besitzt. Wie kann man nun aber die Minimalsuffizienz feststellen? Wir überlegen uns, dass man mithilfe einer Maßzahl  $M(\mathbf{Y})$  den Stichprobenraum in elementfremde Teilmengen derart zerlegen kann, dass alle  $\mathbf{Y}$ , für die  $M(\mathbf{Y})$  den gleichen Wert  $M$  ergibt, derselben Teilmenge angehören. Umgekehrt ist durch eine gegebene Zerlegung auch die Maßzahl definiert. Wir definieren nun eine Zerlegung, von der wir zeigen werden, dass durch sie eine minimal suffiziente Maßzahl gegeben ist.

### Definition 1.7

Es sei  $Y_0 \in \{Y\}$  ein fester Punkt im Stichprobenraum (ein bestimmter Wert einer realisierten Zufallsstichprobe), der die Realisationen einer Zufallsstichprobe  $\mathbf{Y}$  mit Komponenten aus einer Familie  $(P^*(\theta), \theta \in \Omega)$  von Wahrscheinlichkeitsverteilungen enthält. Über die Likelihood-Funktion  $L(\mathbf{Y}, \theta)$  wird durch

$$M(Y_0) = \{Y : L(\mathbf{Y}, \theta) \sim L(Y_0, \theta)\} \quad (1.12)$$

eine Teilmenge in  $\{Y\}$  definiert. Lassen wir  $Y_0$  den ganzen Stichprobenraum  $\{Y\}$  durchlaufen, so wird eine Zerlegung erzeugt. Diese Zerlegung heißt Likelihood-Zerlegung, die ihr entsprechende Maßzahl  $M_L(\mathbf{Y})$  für die  $M_L(\mathbf{Y}) = \text{konst.}$  für alle  $\mathbf{Y} \in M(Y_0)$  und für jedes  $Y_0$  gilt, heißt Likelihood-Maßzahl.

Bevor wir mit dieser Methode minimal suffiziente Maßzahlen für einige Beispiele konstruieren, formulieren wir den

**Satz 1.2**

Die Likelihood-Maßzahl  $M_L(\mathbf{Y})$  ist minimal suffizient bezüglich  $\theta$ .

*Beweis:* Für die Likelihood-Maßzahl  $M_L(\mathbf{Y})$  gilt mit  $Y_1, Y_2 \in \{Y\}$

$$M_L(\mathbf{Y}_1) = M_L(\mathbf{Y}_2)$$

genau dann, wenn  $L(Y_1, \theta) \sim L(Y_2, \theta)$  ist. Damit ist  $L(Y, \theta)$  eine Funktion von  $M_L(\mathbf{Y})$  der Form

$$L(Y, \theta) = a(Y)g^*(M_L(\mathbf{Y}), \theta) \quad (1.13)$$

und nach dem Zerlegungssatz ist  $M_L(\mathbf{Y})$  suffizient bezüglich  $\theta$ . Ist  $M(\mathbf{Y})$  eine beliebige andere bezüglich  $\theta$  suffiziente Maßzahl und gilt für zwei Punkte  $Y_1, Y_2 \in \{Y\}$  die Beziehung  $M(\mathbf{Y}_1) = M(\mathbf{Y}_2)$  sowie  $L(Y_i, \theta) > 0$  mit  $i = 1, 2$ , so folgt ebenfalls aus dem Zerlegungssatz

$$L(Y_1, \theta) = h(Y_1)g(M(\mathbf{Y}_1), \theta) = h(Y_2)g(M(\mathbf{Y}_2), \theta)$$

wegen  $M(\mathbf{Y}_1) = M(\mathbf{Y}_2)$  und  $L(Y_2, \theta) = h(Y_2)g(M(\mathbf{Y}_2), \theta)$  bzw.  $g(M(\mathbf{Y}_2), \theta) = \frac{L(Y_2, \theta)}{h(Y_2)}$ .

Damit wird  $L(Y_1, \theta)$  zu

$$L(Y_1, \theta) = \frac{h(Y_1)}{h(Y_2)}L(Y_2, \theta), h(Y_2) > 0$$

sodass  $L(Y_1, \theta) \sim L(Y_2, \theta)$  ist. Das ist aber gerade die Bedingung dafür, dass  $M(\mathbf{Y}_1) = M(\mathbf{Y}_2)$  ist. Folglich ist  $M_L(\mathbf{Y})$  eine Funktion von  $M(\mathbf{Y})$ , wie  $M(\mathbf{Y})$  auch gewählt wird und damit minimal suffizient.

Wir demonstrieren das Verfahren an zwei Beispielen.

**Beispiel 1.6**

Die Komponenten  $y_i$  einer Zufallsstichprobe  $\mathbf{Y}$  seien nach  $B(N, p)$ ,  $N$  fest,  $0 < p < 1$  binomialverteilt. Es ist eine bezüglich  $p$  minimal suffiziente Maßzahl gesucht. Die Likelihood-Funktion ist

$$L(\mathbf{Y}, p) = \prod_{i=1}^n \binom{N}{y_i} p^{y_i} (1-p)^{(N-y_i)}, \quad y_i = 0, 1, \dots, N$$

Für alle  $\mathbf{Y}_0 = (y_{01}, \dots, y_{0N})^T \in \{Y\}$  mit  $L(\mathbf{Y}_0, p) > 0$  ist

$$\frac{L(\mathbf{Y}, p)}{L(\mathbf{Y}_0, p)} = \frac{\prod_{i=1}^n \binom{N}{y_i}}{\prod_{i=1}^n \binom{N}{y_{0i}}} \left( \frac{p}{1-p} \right)^{\sum_{i=1}^n (y_i - y_{0i})}$$

Damit ist  $M(Y_0)$  auch durch  $M(Y_0) = \{Y : \sum_{i=1}^n y_i = \sum_{i=1}^n y_{0i}\}$  definiert, da gerade dort  $L(Y, p) \sim L(Y_0, p)$  gilt. Folglich ist  $M(Y) = \sum_{i=1}^n y_i$  eine minimal suffiziente Maßzahl.

### Beispiel 1.7

Die Komponenten  $y_i$  einer Zufallsstichprobe  $Y = (y_1, y_2, \dots, y_n)^T$  seien gamma-verteilt. Dann ist für  $y_i > 0$

$$L(Y, a, k) = \frac{a^{nk}}{[\Gamma(k)]^n} e^{-a \sum_{i=1}^n y_i} \prod_{i=1}^n y_i^{k-1}$$

Für alle  $Y_0 = (y_{01}, \dots, y_{0n})^T \in \{Y\}$  mit  $L(Y_0, a, k) > 0$  ist

$$\frac{L(Y, a, k)}{L(Y_0, a, k)} = e^{-a(\sum_{i=1}^n y_i - \sum_{i=1}^n y_{0i})} \frac{\prod_{i=1}^n y_i^{k-1}}{\prod_{i=1}^n y_{0i}^{k-1}}$$

Ist  $a$  vorgegeben, so ist  $\prod_{i=1}^n y_i$  minimal suffizient bezüglich  $k$ . Ist  $k$  bekannt, so ist  $\sum y_i$  minimal suffizient bezüglich  $a$ . Sind  $a$  und  $k$  unbekannte Parameter, so ist  $(\prod_{i=1}^n y_i, \sum_{i=1}^n y_i)$  minimal suffizient bezüglich  $(a, k)$ .

Allgemein gilt:

### Satz 1.3

Ist  $(P^*(\theta), \theta \in \Omega)$  eine  $k$ -parametrische Exponentialfamilie mit der Likelihood-Funktion in kanonischer Form

$$L(y, \theta) = e^{\sum_{i=1}^k \eta_i M_i(y) - A(\eta)} h(y)$$

wobei die Dimension des Parameterraumes gleich  $k$  ist (d. h., die  $\eta_1, \dots, \eta_k$  linear unabhängig sind), dann ist

$$M(Y) = \left( \sum_{i=1}^n M_1(y_i), \dots, \sum_{i=1}^n M_k(y_i) \right)^T$$

minimal suffizient bezüglich  $(P^*(\theta), \theta \in \Omega)$ .

*Beweis:* Die Suffizienz von  $M(Y)$  folgt aus Korollar 1.1 zum Zerlegungssatz, und die Minimalsuffizienz folgt aus der Tatsache, dass  $M(Y)$  die Likelihood-Maßzahl ist, denn es ist genau dann  $L(Y, \theta) \sim L(Y_0, \theta)$ , wenn

$$\sum_{j=1}^k \eta_j \sum_{i=1}^n [M_j(y_i) - M_j(y_{0i})] = 0$$

gilt, und wegen der linearen Unabhängigkeit der  $\eta_i$  ist das nur dann der Fall, wenn  $M(Y) = M(Y_0)$  gilt.



**Beispiel 1.8**

Es sei  $(P^*(\theta), \theta \in \Omega)$  die Familie der zweidimensionalen Normalverteilungen mit der Zufallsvariablen  $\begin{pmatrix} x \\ y \end{pmatrix}$ , dem Erwartungswertvektor  $\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}$  und der Kovarianzmatrix  $\Sigma = \begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}$ . Das ist eine vierparametrische Exponentialfamilie mit den natürlichen Parametern

$$\eta_1 = \frac{\mu_x}{\sigma_x^2}, \quad \eta_2 = \frac{\mu_y}{\sigma_y^2}, \quad \eta_3 = -\frac{1}{2\sigma_x^2}, \quad \eta_4 = -\frac{1}{2\sigma_y^2}$$

und den Faktoren

$$\begin{aligned} M_1 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] &= x, & M_2 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] &= y, \\ M_3 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] &= x^2, & M_4 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] &= y^2, \\ A(\eta) &= \frac{1}{2} \left( \frac{\mu_x^2}{\sigma_x^2} + \frac{\mu_y^2}{\sigma_y^2} \right) \end{aligned}$$

Ist  $\dim(\Omega) = 4$ , so ist

$$M = \left( \sum_{i=1}^n M_{1i}, \quad \sum_{i=1}^n M_{2i}, \quad \sum_{i=1}^n M_{3i}, \quad \sum_{i=1}^n M_{4i} \right)^T$$

minimal suffizient bezüglich  $(P^*(\theta), \theta \in \Omega)$ . Nehmen wir an,  $(\check{P}^*(\theta), \theta \in \Omega) \subseteq (P^*(\theta), \theta \in \Omega)$  sei die Teilfamilie von  $(P^*(\theta), \theta \in \Omega)$ , für die  $\sigma_x^2 = \sigma_y^2 = \sigma^2$  gilt, dann ist  $\dim(\Omega) = 3$ , und  $M$  ist nicht minimal suffizient bezüglich  $\check{P}^*(\theta), \theta \in \Omega$ .

Die natürlichen Parameter von  $\check{P}^*(\theta), \theta \in \Omega$  sind

$$\eta_1 = \frac{\mu_x}{\sigma^2}, \quad \eta_2 = \frac{\mu_y}{\sigma^2}, \quad \eta_3 = -\frac{1}{2\sigma^2}$$

ferner ist  $A(\eta) = \frac{1}{2\sigma^2}(\mu_x^2 + \mu_y^2)$  und die Faktoren der  $\eta_i$  sind

$$\check{M}_1 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] = x, \quad \check{M}_2 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] = y, \quad \check{M}_3 \left[ \begin{pmatrix} x \\ y \end{pmatrix} \right] = x^2 + y^2$$

Bezüglich  $\check{P}^*(\theta), \theta \in \Omega$  ist

$$\check{M} = \left( \sum_{i=1}^n \check{M}_{1i}, \quad \sum_{i=1}^n \check{M}_{2i}, \quad \sum_{i=1}^n \check{M}_{3i} \right)^T$$

minimal suffizient.

Wie in Kapitel 6 am Beispiel des Modells II der Varianzanalyse gezeigt wird, ist das Ergebnis von Satz 1.3 auch in komplizierteren Modellen geeignet, minimal suffiziente Maßzahlen zu finden.

Eine für die Schätztheorie weitere wichtige Eigenschaft ist die Vollständigkeit bzw. beschränkte Vollständigkeit, die wir gemeinsam durch folgende Definition einführen.

**Definition 1.8**

Eine Verteilungsfamilie  $P = (P_\theta, \theta \in \Omega)$  mit der Verteilungsfunktion  $F(y, \theta)$ ,  $\theta \in \Omega$  heißt vollständig, wenn für jede  $P$ -integrierbare Funktion  $h(y)$  der Zufallsvariablen  $y$  aus

$$E[h(y)] = \int h(y) dF(y) = 0 \quad \text{für alle } \theta \in \Omega \quad (1.14)$$

die Beziehung

$$P_\theta[h(y) = 0] = 1 \quad \text{für alle } \theta \in \Omega \quad (1.15)$$

folgt. Folgt (1.15) aus (1.14) nur für beschränkte Funktionen  $h(y)$ , so heißt  $P = (P_\theta, \theta \in \Omega)$  beschränkt vollständig.

Wir wollen ein Beispiel für eine vollständige Verteilungsfamilie betrachten.

**Beispiel 1.9**

Es sei  $P$  die Familie  $\{P_p\}$ ,  $p \in (0, 1)$  der Binomialverteilungen mit der Wahrscheinlichkeitsfunktion

$$p(y, p) = \binom{n}{y} p^y (1-p)^{n-y} = \binom{n}{y} v^y (1-p)^n, \quad 0 < p < 1$$

$$y = 0, 1, \dots, n, \quad v = \frac{p}{1-p}$$

Integrierbarkeit von  $h(y)$  bedeutet Endlichkeit von  $(1-p)^n \sum_{y=0}^n h(y) \binom{n}{y} v^y$  und aus (1.14) folgt

$$\sum_{y=0}^n h(y) \binom{n}{y} v^y = 0 \quad \text{für alle } p \in (0, 1)$$

Das ist ein Polynom  $n$ -ten Grades in  $v$ , das höchstens  $n$  reelle Nullstellen besitzt. Damit diese Gleichung für alle  $v \in \mathbb{R}^+$  erfüllt ist, muss  $\binom{n}{y} h(y)$  für  $y = 0, 1, \dots, n$  verschwinden, und da alle  $\binom{n}{y} > 0$  sind, impliziert das  $P_\theta[h(y) = 0] = 1$  für alle  $p \in (0, 1)$ .

**Satz 1.4**

Eine  $k$ -parametrische Exponentialfamilie der Verteilung der suffizienten Maßzahl ist unter den Voraussetzungen von Satz 1.3 ( $\dim(\Omega) = k$ ) vollständig.

Den Beweis findet man bei Lehmann (1959, S. 132).

**Definition 1.9**

Gegeben sei eine Zufallsstichprobe  $Y = (y_1, y_2, \dots, y_n)^T$ , deren Komponenten einer Verteilung aus der Familie

$$P^* = (P_\theta, \theta \in \Omega)$$

folgen. Eine Maßzahl  $M(Y)$ , deren Verteilung von  $\theta$  unabhängig ist, heißt Hilfsmaßzahl. Ist  $P$  die Familie der durch die Maßzahl  $M(Y)$  aus  $P^*$  induzierten Verteilungen und ist  $P$  vollständig und  $M(Y)$  suffizient bezüglich  $P^*$ , so heißt  $M(Y)$  vollständig suffizient.

**Beispiel 1.10**

Es sei  $P^*$  die Familie der Normalverteilungen  $N(\mu, 1)$  mit Erwartungswert  $\mu = \theta$  und Varianz 1, d. h., es gilt  $\Omega = R^1$ . Das ist eine einparametrische Exponentialfamilie mit  $\dim(\Omega) = 1$ , die nach Satz 1.4 vollständig ist. Ist  $Y = (y_1, y_2, \dots, y_n)^T$  eine Zufallsstichprobe mit Komponenten aus  $P^*$ , so ist  $M_1(Y) = \bar{y}$  nach  $N(\mu, \frac{1}{n})$  verteilt. Die Familie der Verteilungen von  $P^*$  ist folglich auch vollständig. Wegen Satz 1.3 ist  $\bar{y}$  minimal suffizient und damit vollständig suffizient. Die durch  $(n-1)M_2(Y) = \sum y_i^2 - n\bar{y}^2$  induzierte Verteilungsfamilie der  $CQ(n-1)$ -Verteilungen ( $\chi^2$ -Verteilungen mit  $n-1$  Freiheitsgraden) ist von  $\mu$  unabhängig. Folglich ist  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$  bezüglich  $\mu = \theta$  eine Hilfsmaßzahl.

Wir schließen diesen Abschnitt ab mit

**Satz 1.5**

Es sei  $Y$  eine Zufallsstichprobe mit Komponenten aus  $P = (P_\theta, \theta \in \Omega)$  und  $M_1(Y)$  beschränkt vollständig suffizient bezüglich  $P$ . Ist ferner  $M_2(Y)$  eine Maßzahl mit einer von  $\theta$  unabhängigen Verteilung, so sind  $M_1(Y)$  und  $M_2(Y)$  unabhängig.

*Beweis:* Es sei  $\{Y_0\} \subset \{Y\}$  eine Teilmenge des Stichprobenraumes  $\{Y\}$ .  $M_2(Y)$  bildet  $\{Y\}$  auf  $\{M\}$  und  $\{Y_0\}$  auf  $\{M_0\}$  ab. Da die Verteilung von  $M_2(Y)$  von  $\theta$  unabhängig ist, ist  $P[M_2(Y) \in \{M_0\}]$  von  $\theta$  unabhängig. Darüber hinaus ist wegen der Suffizienz von  $M_1(Y)$  bezüglich  $\theta$  auch  $P[M_2(Y) \in \{M_0\} | M_1(Y)]$  von  $\theta$  unabhängig. Wir betrachten die Maßzahl

$$h(M_1(Y)) = P[M_2(Y) \in \{M_0\} | M_1(Y)] - P[M_2(Y) \in \{M_0\}]$$

die von  $M_1(Y)$  abhängt, sodass analog zu (1.14)

$$E_\theta[h(M_1(Y))] = E_\theta[P[M_2(Y) \in \{M_0\} | M_1(Y)] - P[M_2(Y) \in \{M_0\}]] = 0$$

für alle  $\theta \in \Omega$  folgt. Da  $M_1(Y)$  beschränkt vollständig ist, gilt für alle  $\theta \in \Omega$  mit Wahrscheinlichkeit 1 analog zu (1.15)  $P[M_2(Y) \in \{M_0\} | M_1(Y)] - P[M_2(Y) \in \{M_0\}] = 0$  und das bedeutet, dass  $M_1(Y)$  und  $M_2(Y)$  unabhängig sind.

## 1.4

## Der Informationsbegriff in der Statistik

Bei der heuristischen Einführung suffizienter Maßzahlen in Abschn. 1.2 war davon die Rede, dass eine Maßzahl die Information einer Stichprobe weitgehend ausschöpfen sollte. Suffizient wird daher auch mit erschöpfend übersetzt. Was soll aber unter der Information einer Stichprobe eigentlich verstanden werden? Der Informationsbegriff wurde von R.A. Fisher in die Statistik eingeführt, und seine Definition ist auch heute noch von großer Bedeutung. Wir sprechen in diesem Zusammenhang von der Fisher-Information. Ein weiterer Informationsbegriff stammt von Kullback und Leibler (1951), wir wollen aber hier nicht weiter auf diese Definition eingehen. Wir beschränken uns in diesem Abschnitt zunächst auf Verteilungsfamilien

$$P = (P_\theta, \theta \in \Omega), \Omega \subset \mathbb{R}^1$$

mit reellen Parametern  $\theta$ . Mit  $L(y, \theta)$  wird die Likelihood-Funktion ( $Y = y$ ) von  $P$  bezeichnet.

**Definition 1.10**

Es sei  $y$  nach  $P = (P_\theta, \theta \in \Omega), \Omega \subset \mathbb{R}^1$  verteilt. Weiter sei folgende Voraussetzung V1 erfüllt:

1.  $\Omega$  ist ein offenes Intervall.
2. Für jedes  $y \in \{Y\}$  und für jedes  $\theta \in \Omega$  existiert  $\frac{\partial}{\partial \theta} L(y, \theta)$  und ist endlich. Die Menge der Punkte, in denen  $L(y, \theta) = 0$  ist, hängt nicht von  $\theta$  ab.
3. Für jedes  $\theta \in \Omega$  existiert ein  $\varepsilon > 0$  und eine positive  $P_\theta$ -integrierbare Funktion  $k(y, \theta)$  derart, dass für alle  $\theta_0$  in einer  $\varepsilon$ -Umgebung von  $\theta$

$$\left| \frac{L(y, \theta) - L(y, \theta_0)}{\theta - \theta_0} \right| \leq k(y, \theta_0)$$

gilt.

4.  $\frac{\partial}{\partial \theta} L(y, \theta)$  ist quadratisch  $P_\theta$ -integrierbar und es ist für alle  $\theta \in \Omega$

$$0 < E \left\{ \left[ \frac{\partial}{\partial \theta} \ln L(y, \theta) \right]^2 \right\}$$

Dann heißt

$$I(\theta) = E \left\{ \left[ \frac{\partial}{\partial \theta} \ln L(y, \theta) \right]^2 \right\} \quad (1.16)$$

die Fisher-Information der Verteilung  $P_\theta$  bzw. von  $y$ .

Aus der dritten Bedingung von V1 folgt, dass  $P_\theta$ -Integration und Differentiation nach  $\theta$  für  $L(y, \theta)$  vertauscht werden können, und wegen

$$\frac{\partial}{\partial \theta} \ln L(y, \theta) = \frac{\frac{\partial}{\partial \theta} L(y, \theta)}{L(y, \theta)}$$

ist

$$E_{\theta} \left[ \frac{\partial}{\partial \theta} \ln L(y, \theta) \right] = \int_Y \frac{\partial}{\partial \theta} \ln L(y, \theta) L(y, \theta) dy = \frac{\partial}{\partial \theta} \int_Y L(y, \theta) dy = \frac{\partial}{\partial \theta} 1 = 0$$

für alle  $\theta \in \Omega$ . Damit ist

$$I(\theta) = \text{var} \left\{ \frac{\partial}{\partial \theta} \ln L(y, \theta) \right\} \quad (1.17)$$

Es existiere nun auch die zweite Ableitung von  $\ln L(y, \theta)$  nach  $\theta$  für alle  $y$  und  $\theta$ , und es sei  $\int_Y L(y, \theta) dP_{\theta}$  zweifach differenzierbar, wobei man Integration und zweifache Ableitung vertauschen kann. Dann gilt wegen

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} \ln L(y, \theta) &= \frac{L(y, \theta) \frac{\partial^2}{\partial \theta^2} L(y, \theta) - \left( \frac{\partial}{\partial \theta} L(y, \theta) \right)^2}{(L(y, \theta))^2} \\ &= \frac{\frac{\partial^2}{\partial \theta^2} L(y, \theta)}{L(y, \theta)} - \left[ \frac{\frac{\partial}{\partial \theta} L(y, \theta)}{L(y, \theta)} \right]^2 \end{aligned}$$

und

$$0 = \frac{\partial^2}{\partial \theta^2} \int \ln L(y, \theta) dP_{\theta} = \int \frac{\partial^2}{\partial \theta^2} \ln L(y, \theta) dP_{\theta}$$

die Beziehung

$$E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \ln L(y, \theta) \right] = -E_{\theta} \left\{ \left[ \frac{\partial}{\partial \theta} \ln L(y, \theta) \right]^2 \right\} = -I(\theta)$$

und damit

$$I(\theta) = -E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \ln L(y, \theta) \right] \quad (1.18)$$

Wir geben je ein Beispiel für eine diskrete und eine kontinuierliche Verteilung.

### Beispiel 1.11

Es sei  $P$  die Familie der Binomialverteilungen mit gegebenem  $n$  und  $\Omega = (0, 1)$ . Die Likelihood-Funktion ist

$$L(y, p) = \binom{n}{y} p^y (1-p)^{n-y}$$

Die Voraussetzung V1 ist erfüllt, denn das Quadrat von  $\frac{\partial}{\partial p} \ln L(y, p) = \frac{y}{p} - \frac{n-y}{1-p}$  besitzt nach Übergang zu zufälligem  $y$  den endlichen Erwartungswert

$$I(p) = E_p \left\{ \left[ \frac{\partial}{\partial p} \ln L(y, p) \right]^2 \right\} = \sum_{y=0}^n \left( \frac{y}{p} - \frac{n-y}{1-p} \right)^2 \binom{n}{y} p^y (1-p)^{n-y}$$

und daraus folgt

$$I(p) = \frac{n}{p(1-p)}, \quad 0 < p < 1$$

### Beispiel 1.12

Es sei  $P$  die Familie der  $N(\mu, \sigma^2)$ -Verteilungen mit bekanntem  $\sigma^2$ . Wir haben  $\Omega = \mathbb{R}^1$ , und die Likelihood-Funktion hat die Form  $L(y, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$ . Auch für diese Verteilung ist V1 erfüllt. Wir erhalten  $\frac{\partial}{\partial \mu} L(y, \mu) = \frac{1}{\sigma^2}(y - \mu)$  und  $I(\mu) = \frac{1}{\sigma^4} E[(y - \mu)^2] = \frac{1}{\sigma^4} \text{var}(y) = \frac{1}{\sigma^2}$ .

Wir beweisen nun die Additivitätseigenschaft der Fisher-Information.

### Satz 1.6

Existiert für eine Familie  $P$  von Wahrscheinlichkeitsverteilungen mit  $\Omega = \mathbb{R}^1$  die Fisher-Information  $I(\theta) = I_1(\theta)$  und ist  $Y = (y_1, y_2, \dots, y_n)^T$  eine Zufallsstichprobe mit Komponenten  $y_i$  ( $i = 1, \dots, n$ ), die nach  $P_i \in P$  verteilt sind, so ist die Fisher-Information  $I_n(\theta)$  der Verteilung von  $Y$  durch

$$I_n(\theta) = nI_1(\theta) \tag{1.19}$$

gegeben.

*Beweis:* Aus Definition 1.2 folgt, dass die Likelihood-Funktion  $L_n(Y, \theta)$  einer Zufallsstichprobe  $Y$  gleich

$$L_n(Y, \theta) = \prod_{i=1}^n L(y_i, \theta)$$

ist. Damit ist

$$\ln L_n(Y, \theta) = \sum_{i=1}^n \ln L(y_i, \theta)$$

und

$$\frac{\partial}{\partial \theta} \ln L_n(Y, \theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln L(y_i, \theta)$$

Folglich wird wegen (1.17)

$$I_n(\theta) = \text{var} \left\{ \frac{\partial}{\partial \theta} \ln L_n(Y, \theta) \right\} = \sum_{i=1}^n \text{var} \left\{ \frac{\partial}{\partial \theta} \ln L(y_i, \theta) \right\} = nI_1(\theta)$$

**Satz 1.7**

Es sei  $M(Y)$  eine suffiziente Maßzahl bezüglich der Verteilung  $P_\theta \in P$ ,  $\Omega \subseteq R^1$  der Komponenten der Zufallsstichprobe  $Y = (y_1, y_2, \dots, y_n)^T$ . Die Verteilung  $P_\theta$  erfülle die Bedingung V1 von Definition 1.10. Dann existiert die Fisher-Information

$$I_M(\theta) = \left\{ \left[ \frac{\partial}{\partial \theta} \ln L_M(M, \theta) \right]^2 \right\} \quad (1.20)$$

von  $M = M(Y)$ , wobei  $L_M(M, \theta)$  die Likelihood-Funktion von  $M$  ist, und es gilt

$$I_n(\theta) = I_M(\theta) \quad (1.21)$$

*Beweis:* Nach (1.2) bzw. (1.3) ist

$$L(Y, \theta) = h(Y)g(M(Y), \theta)$$

und damit

$$\frac{\partial}{\partial \theta} \ln L(Y, \theta) = \frac{\partial}{\partial \theta} \ln g(M(Y), \theta)$$

da  $h(Y)$  nach Voraussetzung von  $\theta$  unabhängig ist. Wegen Korollar 1.1 von Satz 1.1 gilt für die Likelihood-Funktion  $L_M(M, \theta)$  von  $M$  auch Bedingung V1 von Definition 1.10, und damit existiert  $I_M(\theta)$  von (1.20). Wegen der Äquivalenz von  $L_M(M, \theta)$  und  $L(Y, \theta)$  folgt weiterhin die Behauptung wegen  $\frac{\partial}{\partial \theta} \ln L_M(M, \theta) = \frac{\partial}{\partial \theta} \ln g(M, \theta)$ .

Folglich ist die Fisher-Information einer suffizienten Maßzahl gleich der der Zufallsstichprobe.

Ist  $\theta \in \Omega \subseteq R^p$ , so geben wir folgende

**Definition 1.11**

Es sei  $\mathbf{y}$  nach  $P_\theta \in P$ ,  $\Omega \subseteq R^p$ ,  $\theta = (\theta_1, \dots, \theta_p)^T$  verteilt, und für jede Komponente  $\theta_i$  ( $i = 1, \dots, p$ ) mögen die Bedingungen 2 bis 4 von V1 in Definition 1.10 erfüllt sein.  $\Omega$  sei ein offenes Intervall in  $R^p$ . Ferner existiere der Erwartungswert von  $\frac{\partial}{\partial \theta_i} \ln L(\mathbf{y}, \theta)$   $\frac{\partial}{\partial \theta_j} \ln L(\mathbf{y}, \theta)$  für alle  $\theta$  und alle  $i, j = 1, \dots, p$ . Dann heißt die quadratische Matrix der Ordnung  $p$

$$I(\theta) = (I_{i,j}(\theta)), \quad i, j = 1, \dots, p$$

mit

$$I(\theta) = E \left\{ \frac{\partial}{\partial \theta_i} \ln L(\mathbf{y}, \theta) \frac{\partial}{\partial \theta_j} \ln L(\mathbf{y}, \theta) \right\}$$

die (Fishersche) Informationsmatrix bezüglich  $P_\theta$ .

**Beispiel 1.13**

Die Zufallsvariable  $y$  sei nach  $N(\mu, \sigma^2)$  verteilt mit  $\theta = (\mu, \sigma^2)^T \in R^1 \times R^+ = \Omega$ . Dann ist

$$\ln L(y, \theta) = -\ln \sqrt{2\pi} - \frac{1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (y - \mu)^2$$

und die Voraussetzung von Definition 1.11 ist mit  $\theta_1 = \mu$  und  $\theta_2 = \sigma^2$  erfüllt. Es ist

$$\frac{\partial}{\partial \mu} \ln L(y, \theta) = \frac{y - \mu}{\sigma^2} \quad \text{und} \quad \frac{\partial}{\partial \sigma^2} \ln L(y, \theta) = -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} (y - \mu)^2$$

Wegen  $E[(y - \mu)^2] = \text{var}(y) = \sigma^2$  ist  $I_{11}(\theta) = \frac{1}{\sigma^2}$ , und weil die Schiefe  $\gamma_3 = 0$  ist und wegen  $E(y - \mu) = 0$  ist  $I_{12}(\theta) = I_{21}(\theta) = 0$ . Ferner gilt:

$$\left[ \frac{\partial}{\partial \sigma^2} \ln L(y, \theta) \right]^2 = \frac{1}{4\sigma^4} - \frac{2}{4\sigma^6} (y - \mu)^2 + \frac{1}{4\sigma^8} (y - \mu)^4$$

und daher ist wegen  $E[(y - \mu)^4] = 3\sigma^4$  (denn es ist  $\gamma_2 = 0$ )

$$I_{22} = E \left\{ \left[ \frac{\partial}{\partial \sigma^2} \ln L(y, \theta) \right]^2 \right\} = \frac{1}{\sigma^4} \left[ \frac{1}{4} - \frac{1}{2} + \frac{3}{4} \right] = \frac{1}{2\sigma^4}$$

und wir erhalten

$$I(\theta) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$$

Definieren wir dagegen  $\theta_1 = \mu$  und  $\theta_2 = \sigma$ , so ist  $\frac{\partial}{\partial \sigma} \ln L(y, \theta) = -\frac{1}{\sigma} + \frac{1}{\sigma^2} (y - \mu)^2$ . Während  $I_{11}$ ,  $I_{12}$  und  $I_{21}$  unverändert bleiben, wird

$$I_{22} = E \left\{ \left[ \frac{\partial}{\partial \sigma} \ln L(y, \theta) \right]^2 \right\} = \frac{1}{\sigma^2} [1 - 2 + 3] = \frac{2}{\sigma^2}$$

und damit ist

$$I(\theta) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{pmatrix}$$

Aus diesem Beispiel ersieht man, dass die Fisher-Information nicht invariant gegenüber Parametertransformationen ist. Aus der Kettenregel der Differenzialrechnung folgt allgemein der

**Satz 1.8**

Es sei  $\psi = h(\theta)$  eine in  $\Omega \subseteq R^1$  monotone und bezüglich  $\theta$  differenzierbare Funktion,  $h$  bilde  $\Omega$  auf  $\Pi$  ab. Dann existiert die nach  $\psi$  differenzierbare Umkehrfunktion  $\theta = g(\psi)$ . Unter den Voraussetzungen von Definition 1.10 sei  $I(\theta)$  die Fisher-Information der Verteilung  $P_0 \in \Pi$ . Die Fisher-Information  $I^*(\psi)$  der Verteilung  $P_\psi$  (d. h.,  $P_\theta$  geschrieben mit dem transformierten Parameter) ist dann

$$I^*(\psi) = I(\theta) \left( \frac{d}{d\psi} g(\psi) \right)^2 \quad (1.22)$$



In Beispiel 1.13 war (für festes  $\mu$ )  $\theta = \sigma^2$ ,  $\psi = \sqrt{\theta} = \sigma$  und  $\frac{d\theta}{d\psi} = 2\psi = 2\sigma$ . Nach Satz 1.8 ist also mit  $I(\sigma^2) = \frac{1}{\sigma^2}$

$$I^*(\psi) = I(\sigma^2)4\sigma^2 = \frac{2}{\sigma^2}$$

In Kapitel 2 benötigen wir folgende Ungleichung:

**Satz 1.9** Ungleichung von Rao und Cramér

Es gelte Bedingung VI von Definition 1.10 für die Komponenten der Zufallsstichprobe  $Y$ , deren Likelihood-Funktion  $L(Y, \theta)$  ist. Die Menge  $\{Y_0\} = \{Y \in \{Y\} : L(Y, \theta) = 0\}$  der Punkte des Stichprobenraumes, für die  $L(Y, \theta) = 0$  ist, hänge nicht von  $\theta$  ab. Es sei  $P_\theta \in P = (P_\theta, \theta \in \Omega)$ ,  $\Omega \subseteq R^1$  die Verteilung der Komponenten, und  $M(Y)$  sei eine Maßzahl mit Erwartungswert  $E[M(Y)]$  und Varianz  $\text{var}[M(Y)]$ , die den Stichprobenraum  $\{Y\}$  in  $\Omega$  abbildet. Dann gilt die Rao-Cramér-Ungleichung

$$\text{var}[M(Y)] \geq \frac{\left(\frac{dE[M(Y)]}{d\theta}\right)^2}{nI(\theta)} \quad (1.23)$$

*Beweis:* Mit  $M(Y) = M$  ist  $E = E[M - E(M)] = 0$ , sodass

$$\frac{dE}{d\theta} = - \int_{\{Y\}} \frac{dE[M(Y)]}{d\theta} dP_\theta + \int_{\{Y\}} (M - E(M)) \frac{d}{d\theta} L(Y, \theta) dP_\theta = 0$$

bzw.

$$\frac{dE}{d\theta} = - \frac{dE[M(Y)]}{d\theta} \int_{\{Y\}} dP_\theta + \int_{\{Y\}} (M - E(M)) \frac{d}{d\theta} \ln L(Y, \theta) dY = 0$$

gilt. Daraus folgt

$$E \left\{ (M - E(M)) \frac{d}{d\theta} \ln L(Y, \theta) \right\} = \frac{dE(M)}{d\theta}$$

Wegen der Schwarzischen Ungleichung folgt weiter

$$\begin{aligned} \left\{ \frac{dE(M)}{d\theta} \right\}^2 &= \left\{ E[(M - E(M)) \frac{d}{d\theta} \ln L(Y, \theta)] \right\}^2 \\ &\leq E \left[ (M - E(M))^2 E \left\{ \left[ \frac{d}{d\theta} \ln L(Y, \theta) \right]^2 \right\} \right] \end{aligned}$$

und das ergibt wegen (1.16) den Beweis.

Ist speziell  $E(M) = \theta$ , so hat die Rao-Cramér-Ungleichung die Form

$$\text{var}[M(Y)] \geq \frac{1}{nI(\theta)} \quad (1.24)$$

Wir beweisen noch

**Satz 1.10**

Ist  $\mathbf{y}$  nach einer einparametrischen Exponentialfamilie verteilt, und ist  $g(\theta) = \eta = E(\mathbf{M})$ , so gilt:

$$I^*(\eta) = \frac{1}{\text{var}(\mathbf{M})} \quad (1.25)$$

*Beweis:* Da die Bedingung V1 von Definition 1.10 erfüllt ist, existiert  $I(\eta)$ , und wegen

$$\frac{d}{d\eta} \ln L(Y, \eta) = M(Y) - \frac{d}{d\eta} A(\eta)$$

und nach Satz 1.8 ist  $\text{var}(\mathbf{M}) = I^*(\eta) = I(\theta)[\text{var}(\mathbf{M})]^2$ , und daraus folgt die Behauptung.

Aus der Schwarzschen Ungleichung für zweite Momente für jede Maßzahl  $M(\mathbf{Y})$  mit endlichem zweiten Moment und eine beliebige Funktion  $h(\mathbf{Y}, \theta)$ , deren zweites Moment ebenfalls existiert, folgt:

$$\text{var}(\mathbf{M}) \geq \frac{\text{cov}^2[\mathbf{M}, h(\mathbf{Y}, \theta)]}{\text{var}[h(\mathbf{Y}, \theta)]}$$

**Satz 1.11**

Es sei  $M(\mathbf{Y})$  eine Maßzahl mit dem Erwartungswert  $g(\theta)$  und existierendem zweiten Moment, und  $h_j(\mathbf{Y}, \theta)$ ,  $j = 1, \dots, r$  seien Funktionen, deren zweite Momente existieren. Mit

$$c_j = \text{cov}(M(\mathbf{Y}), h_j), \quad \sigma_{ij} = \text{cov}(h_i, h_j), \quad c^T = (c_1, \dots, c_r)$$

und  $\Sigma = (\sigma_{ij})$ ,  $|\Sigma| \neq 0$  gilt dann stets

$$\text{var}(\mathbf{M}) \geq c^T \Sigma^{-1} c \quad (1.26)$$

*Beweis:* Die Behauptung folgt aus  $\frac{c^T \Sigma^{-1} c}{\text{var}(\mathbf{y})} \leq 1$ .

Mithilfe von (1.26) lässt sich die Rao-Cramér-Ungleichung (1.24) auf den  $p$ -dimensionalen Fall verallgemeinern.

**Satz 1.12**

Die Komponenten einer Zufallsstichprobe  $\mathbf{Y} = (y_1, y_2, \dots, y_n)^T$  seien nach

$$P_\theta \in P = (P_\theta, \theta \in \Omega), \quad \Omega \subseteq R^p, \quad \theta^T = (\theta_1, \dots, \theta_p), \quad p > 1$$

verteilt.  $L(Y, \theta)$  sei die Likelihood-Funktion von  $\mathbf{Y}$ . Ferner seien die Voraussetzungen von Definition 1.10 erfüllt, und die Menge der Punkte in  $\{Y\}$ , für die

$L(Y, \theta) = 0$  ist, hänge nicht von  $\theta$  ab. Es sei  $M(Y)$  eine Maßzahl, deren Erwartungswert  $E[M(Y)] = w(\theta)$  existiert und nach den  $\theta_i$  ableitbar ist. Dann gilt

$$\text{var}[M(Y)] \geq a^T I^{-1} a$$

wobei  $I^{-1}$  die Inverse von  $I(\theta)$  und  $a$  der Vektor der Ableitungen von  $w(\theta)$  nach den  $\theta_i$  ist.

*Beweis:* Da  $I(\theta)$  positiv definit ist und damit  $I^{-1}$  existiert, folgt die Behauptung mit  $h_j = \frac{d}{d\theta_j}$  aus (1.26) und nach Definition 1.11.

## 1.5

### Statistische Entscheidungstheorie

Wir formulieren zunächst das allgemeine statistische Entscheidungsproblem und gehen von einer Menge von Zufallsvariablen  $\{y_t\}$  mit  $t \in R^1$  aus, deren Verteilung  $P_\theta \in P = (P_\theta, \theta \in \Omega)$ ,  $\dim\{\Omega\} = p$  zumindest teilweise unbekannt ist.

Wir beschränken uns hier auf den Fall, dass ausschließlich Aussagen über  $\psi = g(\theta)$  zu machen sind, wobei  $\Omega$  durch  $g$  auf  $Z$  abgebildet wird und  $\dim(Z) = s$  ist.  $Z$  heißt Zustandsraum.

Für die Aussagen über  $\psi$  steht dem Statistiker eine Menge  $\{E\}$  von Entscheidungen zur Verfügung,  $\{E\}$  heißt Entscheidungsraum. Für jedes feste  $t_i$  sei  $Y_{t_i} = (y_{t_i,1}, \dots, y_{t_i,n_i})^T$  eine Zufallsstichprobe vom Umfang  $n_i$ .

Die Gesamtheit der Ergebnisse eines Versuches, anhand dessen eine Entscheidung zu fällen (d. h. aus  $\{E\}$  auszuwählen) ist, sei mit

$$N = \sum_{i=1}^k n_i, A_k = (Y_{t_1}, \dots, Y_{t_k}) \in \prod_{i=1}^k \{Y_{t_i}\} = \{Y_{k,N}\}$$

die Realisation einer Zufallsvariablen  $A_k = (Y_{t_1}, \dots, Y_{t_k})$ .

Nun sei  $d \in D$  eine messbare Abbildung von  $\{Y_{k,N}\}$  auf  $E$ , die jedem  $A_k$  eine Entscheidung  $d(A_k)$  zuordnet,  $d$  heißt Entscheidungsfunktion, und  $D$  ist die Menge der zugelassenen Entscheidungsfunktionen.  $A_k$  wird von der Verteilung von  $A_k$  und den  $k$ -Tupeln  $\mathfrak{S}_k = (t_1, \dots, t_k)$  dem Spektrum des Versuches, und  $\mathfrak{N}_k = (n_1, \dots, n_k)$ , der Belegung des Spektrums, abhängen. Mit

$$\begin{pmatrix} \mathfrak{S}_k \\ \mathfrak{N}_k \end{pmatrix} = \begin{pmatrix} t_1, \dots, t_k \\ n_1, \dots, n_k \end{pmatrix} \in V_n$$

bezeichnen wir den konkreten Versuchsplan, der Element einer Menge  $V$  zugelassener Versuchspläne sei. Außerdem sei eine Verlustfunktion  $L$  als messbare Abbildung von  $E \times Z \times R^1$  in den  $R^m$  gegeben (ihre Festlegung und damit die von  $m$  ist ein außermathematisches Problem), d. h., es ist

$$L = L[d(A_k), \psi, f(M)], \quad d(A_k) \in E, \psi \in Z \quad (1.27)$$

mit der nichtnegativen reellen Funktion  $f(M)$ .

Die Funktion  $L$  gibt den Verlust an, der eintritt, wenn  $d(A_k)$  gewählt wird und  $\psi$  der Wert im transformierten Parameterraum ist,  $f(M)$  entspricht den Kosten für die Realisierung von  $M = (d, \mathfrak{E}_k, \mathfrak{N}_k)$ . Die Aufgabe der Statistik besteht in der Bereitstellung von Methoden zur Auswahl von Tripeln  $M = (d, \mathfrak{E}_k, \mathfrak{N}_k)$ , die ein Risikofunktional  $R$  genanntes Funktional  $R(d, \mathfrak{E}_k, \mathfrak{N}_k, \psi, f(M))$  des zufälligen Verlustes minimieren. Wir werden mit  $d$  entweder eine Entscheidungsfunktion (bei festem  $n$ ) oder eine Folge von Entscheidungsfunktionen gleicher Struktur, deren Elemente sich nur hinsichtlich des Stichprobenumfangs  $n$  unterscheiden, bezeichnen. Wir wollen annehmen, dass

$$R(d, \mathfrak{E}_k, \mathfrak{N}_k, \psi, f(M)) = F(d, \mathfrak{E}_k, \mathfrak{N}_k, \psi) + f(\mathfrak{E}_k, \mathfrak{N}_k) \quad (1.28)$$

gilt, wobei  $f$  nicht von  $d$  abhängt und  $d^*$  die Entscheidungsfunktion (Folge von Entscheidungsfunktionen) ist, für die

$$\min_{d \in D} R(d, \mathfrak{E}_k, \mathfrak{N}_k, \psi) = F(d^*, \mathfrak{E}_k, \mathfrak{N}_k, \psi) \quad (1.29)$$

gilt. Dann kann  $R$  in zwei Schritten minimiert werden. Zunächst bestimmt man  $d^*$  so, dass (1.29) erfüllt ist; im zweiten Schritt bestimmt man  $(\mathfrak{E}_k^*, \mathfrak{N}_k^*)$  so, dass

$$R(d^*, \mathfrak{E}_k^*, \mathfrak{N}_k^*, \psi, f(\mathfrak{E}_k^*, \mathfrak{N}_k^*)) = \min_{\substack{(\mathfrak{E}_k^*) \\ (\mathfrak{N}_k^*) \in V}} R(d^*, \mathfrak{E}_k, \mathfrak{N}_k, \psi, f(\mathfrak{E}_k, \mathfrak{N}_k))$$

gilt.

### Definition 1.12

Ein Tripel  $M^* \in V \times D$  heißt lokal  $R$ -optimal an der Stelle  $\psi_0 \in Z$  bezüglich  $V \times D$ , wenn für alle  $M \in V \times D$

$$R[M^*, \psi_0, f(M^*)] \leq R[M, \psi_0, f(M)]$$

gilt. Ist  $M^*$  für alle  $\psi_0 \in Z$  lokal  $R$ -optimal, so heißt  $M^*$  global  $R$ -optimal.

### Beispiel 1.14

Es sei  $k = 1$  und  $y_{t_1} = y$  nach  $N(\mu, \sigma^2)$  verteilt. Dann ist  $\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \in \Omega = R^1 \times R^+$  und  $A_1 = Y$ . Ferner sei  $\psi = g(\theta) = \mu \in R^1$  und  $d(Y) = \hat{\mu}$  eine statistische Maßzahl mit Realisationen im  $R^1 = E$ . Es sei  $D$  die Klasse der statistischen Maßzahlen mit endlichem zweiten Moment und Realisationen in  $R^1$ . Als Verlustfunktion wählen wir

$$L[\hat{\mu}, \mu, f(T)] = c_1(\hat{\mu} - \mu)^2 + c_2 n K, \quad c_1, c_2, K > 0$$

wobei  $K$  die Kosten einer Messung darstellen. Als Risiko  $R$  wird der erwartete zufällige Verlust

$$R(\hat{\mu}, n, \mu, Kn) = E[c_1(\mu - \hat{\mu})^2 + c_2 n K] = c_2 n K + c_1 [\text{var}(\hat{\mu}) + B(\hat{\mu})^2]$$

gewählt, in dem  $B(\hat{\mu}) = E(\hat{\mu}) - \mu$  ist. In der Klasse  $D$  ist für die Entscheidung  $\hat{\mu} = \psi_0$  zusammen mit  $n = 0$  lokal  $R$ -optimal, für  $(\psi_0, n)$  ist  $R$  gleich 0. Um diesen unbefriedigenden trivialen Fall auszuschließen, kann man  $D$  einschränken. Mit  $D_E \subset D$  bezeichnen wir die Teilmenge in  $D$ , für die  $B(\hat{\mu}) = 0$  ist. Dann wird

$$R(\hat{\mu}, n, \mu, Kn) = c_2 nK + c_1 \text{var}(\hat{\mu}), \quad \hat{\mu} \in D_E$$

und hat die Form (1.28). Wir werden in Kapitel 2 sehen, dass  $\text{var}(\hat{\mu})$  für  $\hat{\mu}^* = \bar{y}$  zum Minimum wird.

Da  $\text{var}(\bar{y}) = \frac{\sigma^2}{n}$  ist, gilt daher im Ergebnis des ersten Schrittes der Minimierung von  $R$

$$\min_{d \in D_E} c_1 \text{var}(\hat{\mu}) = \frac{c_1}{n} \sigma^2$$

und es ist

$$R(\hat{\mu}, n, \mu, Kn) = c_2 nK + \frac{c_1}{n} \sigma^2$$

Leiten wir die rechte Seite nach  $n$  ab und setzen die Ableitung gleich 0, so erhalten wir  $n^* = \sigma \sqrt{\frac{c_1}{Kc_2}}$ , und das hängt ebenso wie  $\bar{y}$  nicht von  $\psi = \mu$  ab, und wegen der Konvexität der Ableitungsfunktion handelt es sich tatsächlich um ein Minimum. Daher ist die in  $Z$  global (jedoch in  $\Omega$  wegen der Abhängigkeit von  $\sigma$  lokal)  $R$ -optimale Lösung des Entscheidungsproblems in  $E \times Z$  gegeben durch

$$M^* = \left( \bar{y}, n^* = \sigma \sqrt{\frac{c_1}{Kc_2}} \right)$$

Wählen wir  $\psi = g(\theta) = \sigma^2 > 0$ , dann ist  $E = R^+$ ,  $k = 1$ ,  $A_1 = Y$  und  $N = n$ . Die Verlustfunktion sei

$$L[d(Y), \sigma^2, f(M)] = c_1(\sigma^2 - d)^2 + c_2 nK, \quad c_i > 0, \quad K > 0$$

Wählen wir als Risiko wieder

$$R(d(Y), n, \sigma^2, Kn) = R = E(L) = c_1 E\{(\sigma^2 - d(Y))^2\} + c_2 nK$$

so ist das wieder von der Form (1.28). Schränken wir uns aus zum vorigen Fall analogen Gründen auf die  $d \in D_E$  ein, für die  $E[d(Y)] = \sigma^2$  gilt, so ist, wie wir in Kapitel 2 sehen werden, der erste Summand von  $R$  für

$$d(Y) = s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

minimal.

Da  $\frac{s^2}{\sigma^2}(n-1)$  nach  $CQ(n-1)$  verteilt ist und damit die Varianz  $2(n-1)$  hat, ist

$$\text{var}(s^2) = \frac{2\sigma^4}{n-1}$$

und es ist nach dem ersten Schritt der Optimierung

$$R(s^2, n, \sigma^2, Kn) = c_1 \frac{2\sigma^4}{n-1} + c_2 nK$$

Das  $R$ -optimale  $n$  ergibt sich zu

$$n^* = 1 + \sigma^2 \sqrt{\frac{2c_1}{Kc_2}}$$

und die lokale  $R$ -optimale Lösung des Entscheidungsproblems ist

$$M^* = \left( s^2, n^* = 1 + \sigma^2 \sqrt{\frac{2c_1}{Kc_2}} \right)$$

Weiterführende theoretische Details und weitere Anwendungsfälle werden in den folgenden Kapiteln bei der Wahl des minimalen Stichprobenumfangs behandelt. Wir wollen davon ausgehen, dass  $d$  bei festem  $\mathfrak{S}_k$  und  $\mathfrak{N}_k$  bezüglich einer bestimmten Risikofunktion  $R$ -optimal zu wählen ist. Bezüglich der optimalen Wahl von  $(\mathfrak{S}_k)$  verweisen wir auf die Kapitel 8 und 9 zur Regressionsanalyse. Wir schreiben daher

$$R(d, \psi) = E\{L[d(Y), \psi]\} = r(d, \tau) \quad (1.30)$$

Um trivial lokal  $R$ -optimale Entscheidungsfunktionen  $d$  zu vermeiden, wurde in Beispiel 1.14 eine Einschränkung auf eine Teilklasse  $D_E \subseteq D$  vorgenommen. Hier sollen zwei weitere allgemeine Vorgehensweisen zur Überwindung solcher Probleme vorgestellt werden.

### Definition 1.13

Es sei  $\theta$  eine Zufallsvariable mit Realisationen  $\theta \in \Omega$  mit der Wahrscheinlichkeitsverteilung  $P_\tau$ ,  $\tau \in \mathfrak{T}$ . Bezüglich  $P_\tau$  möge der Erwartungswert von (1.30)

$$\int_{\Omega} R(d, \psi) d\Pi_\tau = r(d, \tau) \quad (1.31)$$

der Bayessches Risiko bezüglich der a-priori-Verteilung  $P_\tau$  genannt wird, existieren.

Eine Entscheidungsfunktion  $d_0(Y)$ , die

$$r(d_0, \tau) = \min_{d \in D} [r(d, \tau)]$$

erfüllt, heißt Bayessche Entscheidungsfunktion bezüglich der a-priori-Verteilung  $P_\tau$ .

**Definition 1.14**

Eine Entscheidungsfunktion  $d_0 \in D$  heißt Minimax-Entscheidungsfunktion, wenn

$$\max_{\theta \in \Omega} R(d_0, \psi) = \min_{d \in D} \max_{\theta \in \Omega} R(d, \psi) \quad (1.32)$$

gilt.

**Definition 1.15**

Es seien  $d_1, d_2 \in D$  Entscheidungsfunktionen für ein bestimmtes Entscheidungsproblem mit der Risikofunktion  $R(d, \psi)$  mit  $\psi = g(\theta)$ ,  $\theta \in \Omega$ . Dann heißt  $d_1$  nicht schlechter als  $d_2$ , wenn  $R(d_1, \psi) \leq R(d_2, \psi)$  für alle  $\theta \in \Omega$  gilt,  $d_1$  heißt besser als  $d_2$ , wenn neben  $R(d_1, \psi) \leq R(d_2, \psi)$  für alle  $\theta \in \Omega$  für wenigstens ein  $\theta^* \in \Omega$  die Ungleichung  $R(d_1, \psi^*) < R(d_2, \psi^*)$  mit  $\psi^* = g(\theta^*)$  gilt. Eine Entscheidungsfunktion  $d$  heißt zulässig in  $D$ , wenn es in  $D$  keine Entscheidungsfunktion gibt, die besser als  $d$  ist. Ist eine Entscheidungsfunktion nicht zulässig, so heißt sie unzulässig.

Eine weitere Darstellung der Entscheidungstheorie ist hier nicht erforderlich. Wir werden in Kapitel 2 die Theorie der Punktschätzungen behandeln, dort ist  $d(Y) = S(Y)$  eine Entscheidungsfunktion. In der Testtheorie in Kapitel 3 ist  $d(Y)$  die Wahrscheinlichkeit für die Ablehnung einer Nullhypothese und in der Konfidenzschätzung ein Bereich in  $\Omega$ , der den Wert  $\theta$  der Verteilung  $P_\theta$  mit einer vorgegebenen Wahrscheinlichkeit überdeckt. Auswahlregeln und multiple Vergleichsverfahren sind andere Spezialfälle von Entscheidungsfunktionen.

**1.6****Übungsaufgaben****Aufgabe 1.1**

Um das Durchschnittseinkommen der Bewohner einer Großstadt zu schätzen, wird das Einkommen der Besitzer jedes 20. Privatanschlusses in einem Telefonbuch ermittelt.

Handelt es sich bei dieser Stichprobe um eine Zufallsstichprobe der Bevölkerung der Stadt?

**Aufgabe 1.2**

Aus einer Grundgesamtheit mit den Elementen 1, 2, 3 kann man mit Zurücklegen  $3^4 = 81$  verschiedene Stichproben vom Umfang  $n = 4$  auswählen. Man schreibe alle möglichen Stichproben auf, berechne  $\bar{y}$  und  $s^2$  und stelle die Häufigkeitsverteilung von  $\bar{y}$  und  $s^2$  als Streifendiagramm dar.

**Aufgabe 1.3**

Man beweise, dass die jeweilige Maßzahl  $M(Y)$  suffizient bezüglich  $\theta$  ist, wobei  $Y = (y_1, y_2, \dots, y_n)^T$   $n \geq 1$  eine Zufallsstichprobe aus einer Grundgesamtheit mit der Verteilung  $P_\theta$  mit  $\theta \in \Omega$  ist, indem man die bedingte Verteilung von  $Y$  bei gegebenem  $M(Y)$  bildet.

- $M(Y) = \sum_{i=1}^n y_i$  und  $P_\theta$  ist die Poisson-Verteilung mit dem Parameter  $\theta \in \Omega \subset \mathbb{R}^+$ .
- $M(Y) = (y_{(1)}, y_{(n)})^T$  und  $P_\theta$  ist die Gleichverteilung im Intervall  $(\theta, \theta + 1)$  mit  $\theta \in \Omega \subset \mathbb{R}^1$ .
- $M(Y) = y_{(n)}$  und  $P_\theta$  ist die Gleichverteilung im Intervall  $(0, \theta)$  mit  $\theta \in \Omega = \mathbb{R}^+$ .
- $M(Y) = \sum_{i=1}^n y_i$  und  $P_\theta$  ist die Exponentialverteilung mit dem Parameter  $\theta \in \Omega = \mathbb{R}^+$ .

**Aufgabe 1.4**

Es sei  $Y = (y_1, y_2, \dots, y_n)^T$   $n \geq 1$  eine Zufallsstichprobe aus einer Grundgesamtheit mit der Verteilung  $P_\theta$ ,  $\theta \in \Omega$ . Man bestimme mithilfe des Korollars 1.1 zum Zerlegungssatz eine suffiziente Maßzahl bezüglich  $\theta$ , wenn  $P_\theta$ ,  $\theta \in \Omega$  die Dichtefunktion

- $f(y, \theta) = \theta y^{\theta-1}$ ,  $0 < y < 1$ ;  $\theta \in \Omega = \mathbb{R}^+$
- der Weibull-Verteilung

$$f(y, \theta) = \theta a (\theta y)^{a-1} e^{-(\theta y)^a}, \quad y \geq 0, \quad \theta \in \Omega = \mathbb{R}^+, a > 0 \text{ bekannt}$$

- der Pareto-Verteilung

$$f(y, \theta) = \frac{\theta a^\theta}{y^{\theta+1}}, \quad y > a > 0, \quad \theta \in \Omega = \mathbb{R}^+ \text{ bekannt}$$

besitzt.

**Aufgabe 1.5**

Man bestimme eine minimal suffiziente Maßzahl  $M(Y)$  für den Parameter  $\theta$ , wenn  $Y = (y_1, y_2, \dots, y_n)^T$   $n \geq 1$  eine Zufallsstichprobe aus einer Grundgesamtheit mit der folgenden Verteilung  $P_\theta$  ist:

- geometrische Verteilung mit der Wahrscheinlichkeitsfunktion

$$p(y, p) = p(1-p)^{y-1}, \quad y = 1, 2, \dots, 0 < p < 1$$

- hypergeometrische Verteilung mit der Wahrscheinlichkeitsfunktion

$$p(y, M, N, n) = \frac{\binom{M}{y} \binom{N-M}{n-y}}{\binom{N}{n}}, \quad n \in \{1, \dots, N\}$$



c) negative Binomialverteilung mit der Wahrscheinlichkeitsfunktion

$$p(y, p, r) = \binom{y-1}{r-1} p^r (1-p)^{y-r}, \quad 0 < p < 1, y \geq r \text{ ganz,}$$

$$r \in \{0, 1, \dots\}$$

und i)  $\theta = p$  und  $b$  bekannt; ii)  $\theta^T = (p, r)$ .

d) Betaverteilung mit der Dichtefunktion

$$f(y, \theta) = \frac{1}{B(a, b)} y^{a-1} (1-y)^{b-1}, \quad 0 < y < 1, \quad 0 < a, b < \infty$$

und i)  $\theta = a$  aber  $b$  bekannt; ii)  $\theta = b$  aber  $a$  bekannt.

### Aufgabe 1.6

Man beweise, dass die folgenden Verteilungsfamilien  $\{P_\theta, \theta \in \Omega\}$  vollständig sind:

- $P_\theta$  ist die Poisson-Verteilung mit dem Parameter  $\theta \in \Omega = \mathbb{R}^+$ .
- $P_\theta$  ist die Gleichverteilung im Intervall  $(0, \theta)$ ,  $\theta \in \Omega = \mathbb{R}^+$ .

### Aufgabe 1.7

Es sei  $Y = (y_1, y_2, \dots, y_n)^T$   $n \geq 1$  eine Zufallsstichprobe, deren Komponenten im Intervall  $(0, \theta)$ ,  $\theta \in \Omega = \mathbb{R}^+$  gleichverteilt sind. Man zeige, dass  $M(Y) = \mathbf{y}_{(n)}$  vollständig suffizient ist.

### Aufgabe 1.8

Es besitze  $y$  die diskrete Verteilung  $P_\theta$  mit der Wahrscheinlichkeitsfunktion

$$p(y, \theta) = P(y = y) = \begin{cases} \theta & \text{für } y = -1 \\ (1-\theta)^2 \theta^y & \text{für } y = 0, 1, 2 \end{cases}$$

Man zeige, dass die entsprechende Verteilungsfamilie mit  $\theta \in (0, 1)$  beschränkt vollständig, aber nicht vollständig ist.

### Aufgabe 1.9

Gegeben sei eine einparametrische Exponentialfamilie mit der Dichte- oder Wahrscheinlichkeitsfunktion

$$f(y, \theta) = h(y) e^{\{\eta(\theta)M(y) - B(\theta)\}}, \quad \theta \in \Omega$$

- Man drücke die Fisher-Information dieser Verteilung durch die Funktionen  $\eta(\theta)$  und  $B(\theta)$  aus.
- Man benutze das Ergebnis von a) zur Berechnung von  $I(\theta)$  für die
  - Binomialverteilung mit dem Parameter  $\theta = p$ ,
  - Poisson-Verteilung mit dem Parameter  $\theta = \lambda$ ,
  - Exponentialverteilung mit dem Parameter  $\theta$ ,
  - Normalverteilung  $N(\mu, \sigma^2)$  mit  $\theta = \sigma, \mu$  fest.

**Aufgabe 1.10**

Es seien die Voraussetzungen aus Definition 1.11 erfüllt. Außerdem mögen die zweiten Ableitungen  $\frac{\partial^2}{\partial\theta_i\partial\theta_j}L(y, \theta)$  für alle  $i, j = 1, \dots, p$  und  $y \in \{Y\}$  und ihre Erwartungswerte für zufälliges  $y$  existieren, und es sei  $\int_{\{Y\}} L(y, \theta) dy$  zweimal differenzierbar, wobei Integration und Differentiation vertauscht werden können. Man beweise, dass dann die Elemente der Informationsmatrix in Definition 1.11 die Gestalt

$$\text{a) } I_{i,j}(\theta) = \text{cov} \left[ \frac{\partial}{\partial\theta_i} \ln L(y, \theta), \frac{\partial}{\partial\theta_j} \ln L(y, \theta) \right]$$

$$\text{b) } I_{i,j}(\theta) = -E \left[ \frac{\partial^2}{\partial\theta_i\partial\theta_j} \ln L(y, \theta) \right]$$

besitzen.

**Aufgabe 1.11**

Es sei  $Y = (y_1, y_2, \dots, y_n)^\top$  eine Zufallsstichprobe aus einer Grundgesamtheit mit der Verteilung  $P_\theta$ ,  $\theta \in \Omega$  und  $M(Y)$  eine gegebene Maßzahl. Man berechne  $E[M(Y)]$ ,  $\text{var}[M(Y)]$ , die Fisher-Information  $I(\theta)$  der Verteilung und die Rao-Cramér-Schranke für  $\text{var}[M(Y)]$ .

Gilt in der Rao-Cramér-Ungleichung das Gleichheitszeichen, wenn

a)  $P_\theta$  die Poisson-Verteilung mit dem Parameter  $\theta \in R^+$  und

$$M(Y) = \begin{cases} 1 & \text{für } y = 0 \\ 0 & \text{sonst} \end{cases}$$

ist (hier ist  $n = 1$ , d. h.  $y = Y$ );

b)  $P_\theta$  die Poisson-Verteilung mit dem Parameter  $\theta \in R^+$  und  $M(Y) = (1 - \frac{1}{n})^{n\bar{y}}$ . (Verallgemeinerung von a) auf den Fall  $n > 1$ ) ist;

c)  $f(y, \theta) = \theta y^{\theta-1}$ ,  $0 < y < 1$ ,  $\theta \in R^+$  die Dichtefunktion von  $P_\theta$  und  $M(Y) = -\frac{1}{n} \sum_{i=1}^n \ln y_i$  ist?

**Aufgabe 1.12**

In einem Gebiet soll nach Öl gebohrt werden. Der Besitzer der Bohrrechte muss sich auf eine Strategie aus  $E = \{E_1, E_2, E_3\}$  festlegen. Dabei bedeute:  $E_1$  – Bohrung wird selbst durchgeführt,  $E_2$  – die Bohrrechte werden verkauft,  $E_3$  – ein Teil der Bohrrechte wird veräußert.

Es ist jedoch nicht bekannt, ob in dem Gebiet tatsächlich Öl vorkommt.

Es sei  $\Omega = \{\theta_1, \theta_2\}$  wobei  $\theta = \theta_1$  – Öl ist dort vorhanden,  $\theta = \theta_2$  – Öl ist dort nicht vorhanden, bedeuten soll. Die Verlustfunktion  $L(d, \theta)$  hat für die Entscheidungen  $d = E_i$ ,  $i = 1, 2, 3$  und  $\theta = \theta_j$ ,  $j = 1, 2$ , die Form

	$E_1$	$E_2$	$E_3$
$\theta_1$	0	10	5
$\theta_2$	12	1	6

Die Entscheidung wird aufgrund von Gutachten über die geologischen Verhältnisse in dem Gebiet getroffen: Das Ergebnis der Gutachten sei durch  $y \in \{0, 1\}$  gekennzeichnet.

Die Wahrscheinlichkeitsfunktion – in Abhängigkeit von  $\theta$  – der Zufallsvariablen  $y$  sei  $p_\theta(y)$  mit den Werten

	$y = 0$	$y = 1$
$\theta_1$	0,3	0,7
$\theta_2$	0,6	0,4

$y$  gibt also die aus dem „Zufallsexperiment“ der geologischen Gutachten erhaltene Information über das Vorhandensein ( $y = 1$ ) oder Fehlen ( $y = 0$ ) von Ölvorkommen in dem Gebiet an. Die Menge  $D$  der Entscheidungsfunktionen  $d(y)$  enthalte alle nur möglichen  $3^2$  diskreten Funktionen:

	1	2	3	4	5	6	7	8	9
$d_i(0)$	$E_1$	$E_1$	$E_1$	$E_2$	$E_2$	$E_2$	$E_3$	$E_3$	$E_3$
$d_i(1)$	$E_1$	$E_2$	$E_3$	$E_1$	$E_2$	$E_3$	$E_1$	$E_2$	$E_3$

- Man bestimme das Risiko  $R(d(y), \theta) = E_\theta[L\{d(y), \theta\}]$  für alle obigen 18 Fälle.
- Man ermittle die Minimax-Entscheidungsfunktion.
- Nach Meinung von Experten der Bohrtechnik ist die Wahrscheinlichkeit, bei der Niederbringung einer Bohrung in diesem Gebiet auf Öl zu stoßen, gleich 0,2. Dann kann  $\theta$  als Zufallsvariable mit der Wahrscheinlichkeitsfunktion

$\theta$	$\theta_1$	$\theta_2$
$\pi(\theta)$	0,2	0,8

betrachtet werden. Man bestimme für jede Entscheidungsfunktion das Bayessche Risiko  $r(d_i, \pi)$  und anschließend die Bayessche Entscheidungsfunktion.

### Aufgabe 1.13

Es sollen die Behandlungsstrategien beim Einsatz zweier Medikamente  $M_1$  und  $M_2$  beurteilt werden. Drei derartige Strategien stehen zur Verfügung:  $E_1$  – Behandlung mit dem blutdruckerhöhenden Medikament  $M_1$ ;  $E_2$  – Behandlung ohne Medikamente;  $E_3$  – Behandlung mit dem blutdrucksenkenden Medikament  $M_2$ ;  $\theta$  charakterisiert den (geeignet transformierten) Blutdruck eines Patienten:  $\theta < 0$  zu niedriger Blutdruck,  $\theta = 0$  Blutdruck normal,  $\theta > 0$  zu hoher Blutdruck.

Die Verlustfunktion ist folgendermaßen definiert:

	$E_1$	$E_2$	$E_3$
$\theta < 0$	0	$c$	$b + c$
$\theta = 0$	$b$	0	$b$
$\theta > 0$	$b + c$	$c$	0

Bei einem Patienten wird der Blutdruck gemessen. Die Messung  $y$  sei nach  $N(\theta, 1)$  verteilt und wird  $n$ -mal unabhängig voneinander durchgeführt:  $Y = (y_1, y_2, \dots, y_n)^T$ , aufgrund dieser Stichprobe wird die Entscheidungsfunktion

$$d_{r,s} = \begin{cases} E_1, & \text{falls } \bar{y} < r \\ E_2, & \text{falls } r \leq \bar{y} \leq s \\ E_3, & \text{falls } \bar{y} > s \end{cases}$$

definiert.

- Man bestimme die Risikofunktion  $R(d_{r,s}(\bar{y}), \theta) = E\{L[d_{r,s}(\bar{y}), \theta]\}$ .
- Man skizziere die Risikofunktion im Fall  $b = c = 1$ ,  $n = 1$  für i)  $r = -s = -1$ ;  
ii)  $r = -\frac{1}{2}s = -1$ .

Für welche Werte von  $\theta$  ist die Entscheidungsfunktion  $d_{-1,1}(y)$  der Funktion  $d_{-1,2}(y)$  vorzuziehen?

## Literatur

- Bahadur, R.R. (1955) Statistics and subfields. *Ann. Math. Stat.*, **26**, 490–497.
- Blackwell, D. (1947) Conditional expectations and unbiased sequential estimation. *Ann. Math. Stat.*, **18**, 105–110.
- Cochran, W.G. und Boing, W. (1972) *Stichprobenverfahren*, De Gruyter, Berlin, New York.
- Fisher, R.A. (1925) *Statistical Methods for Research Workers*, Oliver & Boyd, Edinburgh.
- Halmos, P.R. und Savage, L.J. (1949) Application of the Radon-Nykodin theorem to the theory of sufficient statistics. *Ann. Math. Stat.*, **20**, 225–241.
- Kauermann, G. und Küchenhoff, H. (2011) *Stichproben: Methoden und praktische Umsetzung mit R*, Springer, Heidelberg.
- Kullback, S. und Leibler, R.A. (1951) On information and sufficiency. *Ann. Math. Stat.*, **22**, 79–86.
- Lehmann, E.L. und Romano, J.P. (2008) *Testing Statistical Hypothesis*, Springer, Heidelberg.
- Lehmann, E.L. und Scheffé, H. (1950) Completeness, similar regions and unbiased estimation. *Sankhya*, **10**, 305–340.
- Quatember, A. (2014) *Datenqualität in Stichprobenerhebungen*, Springer, Berlin.
- Rao, C. R. (1945) Information and accuracy attainable in estimation of statistical parameters. *Bull. Calc. Math. Soc.*, **37** (3), 81–91.
- Rasch, D. (1995) *Mathematische Statistik*, Joh. Ambrosius Barth, Berlin, Heidelberg
- Rasch, D., Tiku, M.L. und Sumpff, D. (Hrsg.) (1994) *Elsevier's Dictionary of Biometry*, Elsevier, Amsterdam, London, New York.
- Rasch, D., Herrendörfer, G., Bock, J., Victor, N. und Guiard, V. (Hrsg.) (2008) *Verfahrensbibliothek Versuchsplanung und -aus-*

wertung, 2. verbesserte Auflage in einem Band mit CD, R. Oldenbourg, München, Wien  
(frühere Auflagen mit den Herausgebern Rasch, Herrendörfer, Bock, Busch (1978, 1981), Deutscher Landwirtschaftsverlag

Berlin und (1995, 1996) Oldenbourg, München Wien).

Stigler, S.M. (1986, 1990) *The History of Statistics: The Measurement of Uncertainty Before 1900*, Harvard University Press, Cambridge.

