

Advanced Measure Theory
and
Mathematical Statistics

Lutz Dümbgen

University of Bern
Spring Semester 2022

October 20, 2022

Bibliography

- [1] E.M. ALFSEN (1963). A simplified constructive proof of the existence and uniqueness of Haar measure. *Mathematica Scandinavica* **12**, pp. 106-116.
- [2] G. AILAM (1968). On probability properties of random sets and the asymptotic behavior of empirical distribution functions. *Journal of Applied Probability* **5(1)**, pp. 196-202.
- [3] H. BAUER (2001). *Measure and integration theory*. Walter de Gruyter & Co., Berlin.
- [4] C. CZADO and T. SCHMIDT (2011). *Mathematische Statistik*. Springer.
- [5] L. DÜMBGEN (2021). *Optimization Methods*. Lecture notes, University of Bern.
- [6] A. HAAR (1933). Der Massbegriff in der Theorie der Kontinuierlichen Gruppen. *Annals of Mathematics* **34(1)**, pp. 147-169.
- [7] P.R. HALMOS (1946). The theory of unbiased estimation. *Annals of Mathematical Statistics* **17**, pp. 34-43.
- [8] W. Hoeffding (1948). A class of statistics with asymptotically normal distribution. *Annals of Mathematical Statistics* **19(3)**, pp. 293-325.
- [9] M.G. KENDALL (1938). A new measure of rank correlation. *Biometrika* **30(1-2)**, pp. 81-93.
- [10] D.W. MÜLLER (1986). *Mathematische Statistik*. Lecture notes, University of Heidelberg.

Acknowledgements. The second part of these lecture notes is an updated and translated version of my personal notes from a lecture “Mathematische Statistik” by my former PhD advisor Prof. Dietrich Werner Müller at the University of Heidelberg in the summer of 1986. His intriguing lectures stimulated my own interest in statistics as a mathematical as well as applied discipline. Constructive comments by Stefan Binder, Alessandro Corrent, Claire Descombes, Alexander Henzi, Xinwei Li, Jin M. Löffler, Christof Mahnig, Alexandre Mösching, Géraldine Oppliger, Levi Ryffel, Sara Salvador, Nadine Seitlinger, Christof Strähl and Nikita van Gils on the current notes are gratefully acknowledged.

Contents

I	Advanced Measure Theory	11
1	Abstract Integrals	13
1.1	Lattices and Stone Lattices	13
1.2	Abstract and Usual Integrals	15
1.3	Representations of Dual Spaces	20
2	Signed Measures	25
2.1	The Hahn–Jordan Decomposition	25
2.2	Radon–Nikodym Derivatives	28
2.2.1	Finite measures	30
2.2.2	Absolute continuity and σ -finite measures	32
2.3	Another Riesz Representation	34
3	Conditional Expectations	37
3.1	Conditional expectations with respect to a sub- σ -field	37
3.2	Conditional expectations as orthogonal projections	40
3.3	Conditional expectations given another random variable	41
4	Stochastic Kernels	47
4.1	Stochastic Kernels and Fubini’s Theorem	47
4.2	Decomposing Measures on Product Spaces	51
4.3	Conditional Expectations and Distributions	55
5	Haar Measure	59
5.1	Locally Compact Topological Groups	59
5.2	Left- and Right-Invariant Measures	62
5.3	Some Explicit Constructions	69

II	Mathematical Statistics	73
6	Measurement Series and Estimators of Location	75
6.1	Statistical Experiments and Point Estimators	75
6.2	Estimators of Location	76
6.3	Constructing an Optimal Equivariant Estimator	77
6.4	Beyond Equivariance: Admissibility	85
6.5	Location Functionals and Gross Error Models	90
7	Statistical Tests	95
7.1	The Neyman–Pearson Lemma	97
7.2	Monotone Density Ratios	100
7.3	Stochastic Order, P-Values, Confidence Bounds	105
7.4	The Generalized Neyman–Pearson Lemma	109
7.5	Tests of Two-Sided Hypotheses	111
7.5.1	One-parameter exponential families (with natural parametrization)	111
7.5.2	Two-sided hypotheses, version 1	113
7.5.3	Two-sided hypotheses, version 2	117
7.5.4	Summary and some first applications	120
7.6	Tests and Confidence regions	124
8	Decision Problems and Procedures, Sufficiency and Completeness	129
8.1	Decision Problems and Procedures	129
8.2	Some Optimality Concepts and Results	132
8.3	Informativity and Sufficiency	136
8.3.1	Informativity	136
8.3.2	Sufficiency	138
8.4	Complete Statistical Experiments	144
8.5	U -Statistics	148
9	Exponential Families	159
9.1	Definitions and Basic Properties	159
9.2	Nuisance Parameters	161

10 Some Asymptotics	175
10.1 Testing, Total Variation and Hellinger Distances	175
10.2 Asymptotics for Repeated Binary Experiments	180
10.3 Fisher Information	188
A Auxiliary Results	195
A.1 Some Basics from Measure Theory	195
A.2 Two Compactness Properties of Statistical Tests	199
A.3 Uniqueness of Moment-Generating Functions	201
A.4 Hoeffding's Decomposition	202
A.5 Weak Law of Large Numbers and Central Limit Theorem	206

Introduction

Statistics is the art of analysing data and dealing with non-avoidable errors and uncertainties in a concise way. In introductory and many advanced Statistics courses, various procedures such as point estimators, statistical tests and confidence regions are introduced for different settings, but often they seem a bit ad hoc. The purpose of Mathematical Statistics is to present these procedures in a coherent framework and to clarify which procedures are optimal for a given task. This includes the question of how to quantify the quality of a statistical procedure.

An indispensable tool for mathematical statistics is measure theory, including Radon-Nikodym derivatives, conditional expectations, conditional distributions and Markov kernels. Hence, the first part of this course is devoted to these aspects of measure theory.

Part I

Advanced Measure Theory

Chapter 1

Abstract Integrals

We tacitly assume that the reader is familiar with standard measure and integration theory. Some basic notions and results are listed in Section A.1.

In standard measure theory, one starts from a measurable space (Ω, \mathcal{A}) , consisting of a set Ω and a σ -field \mathcal{A} over Ω . Then one constructs or defines a measure μ on \mathcal{A} , and this leads to the integral $\int f d\mu$ of measurable functions $f : \Omega \rightarrow [0, \infty]$ or $f : \Omega \rightarrow \mathbb{R}$. If we restrict our attention to the set $\mathcal{L}^1(\mu)$ of real-valued, measurable functions f on Ω such that $\int |f| d\mu < \infty$, then the integral $\int f d\mu$ has the following essential properties:

Linearity: The set $\mathcal{L}^1(\mu)$ is a real vector space¹, and $f \mapsto \int f d\mu$ is linear in $f \in \mathcal{L}^1(\mu)$.

Positivity: If $f \in \mathcal{L}^1(\mu)$ is nonnegative, then $\int f d\mu \geq 0$.

Montone convergence: If $(f_n)_n$ is a sequence of nonnegative functions in $\mathcal{L}^1(\mu)$ which is pointwise increasing with limit $f \in \mathcal{L}^1(\mu)$, then $\int f_n d\mu \rightarrow \int f d\mu$ as $n \rightarrow \infty$.

In the present section we follow the reverse route. Throughout let \mathcal{F} be a real vector space of real-valued functions on a set Ω . Suppose that $J : \mathcal{F} \rightarrow \mathbb{R}$ is a linear and positive functional satisfying monotone convergence. Does there exist a σ -field \mathcal{A} over Ω and a measure μ on \mathcal{A} such that $\mathcal{F} \subset \mathcal{L}^1(\mu)$ and $J(f) = \int f d\mu$ for all $f \in \mathcal{F}$?

1.1 Lattices and Stone Lattices

For two functions $f, g : \Omega \rightarrow [-\infty, \infty]$, the inequalities $f \leq g$ or $f \geq g$ are always meant pointwise. Their pointwise maximum and minimum are denoted with $f \vee g$ and $f \wedge g$, respectively, that is,

$$f \vee g(\omega) := \max\{f(\omega), g(\omega)\} \quad \text{and} \quad f \wedge g(\omega) := \min\{f(\omega), g(\omega)\}.$$

Similarly, the pointwise maximum and minimum of f and a real constant c is denoted with $f \vee c$ and $f \wedge c$, respectively. Specifically, we write

$$f^+ := f \vee 0 \quad \text{and} \quad f^- := (-f) \vee 0 = -(f \wedge 0),$$

so $f = f^+ - f^-$ and $|f| = f^+ + f^-$.

¹with the pointwise addition of functions and the pointwise multiplication of functions with scalars

Definition 1.1 (Lattice and Stone lattice). The linear function space \mathcal{F} is called a *lattice* if $|f| \in \mathcal{F}$ for any $f \in \mathcal{F}$.

A lattice \mathcal{F} is called a *Stone lattice* if $f \wedge 1 \in \mathcal{F}$ for arbitrary $f \in \mathcal{F}$.

Example 1.2. Let (Ω, d) be a metric space. Then the following sets of functions $f : \Omega \rightarrow \mathbb{R}$ are Stone lattices:

$$\begin{aligned}\mathcal{C}(\Omega) &= \{f : \text{continuous}\}, \\ \mathcal{C}_b(\Omega) &= \{f : \text{continuous and bounded}\}, \\ \mathcal{C}_{\text{Lip}}(\Omega) &= \{f : \text{Lipschitz-continuous}\}, \\ \mathcal{C}_{\text{Lip},b}(\Omega) &= \{f : \text{Lipschitz-continuous and bounded}\}.\end{aligned}$$

Exercise 1.3 (Lattices). Show that the following three properties of a linear function space \mathcal{F} are equivalent:

- (i) $|f| \in \mathcal{F}$ for arbitrary $f \in \mathcal{F}$.
- (ii) $f^+ \in \mathcal{F}$ for arbitrary $f \in \mathcal{F}$.
- (iii) $f \vee g, f \wedge g \in \mathcal{F}$ for arbitrary $f, g \in \mathcal{F}$.

From now on, let \mathcal{F} be a Stone lattice. The set of nonnegative functions in \mathcal{F} is denoted with \mathcal{F}^+ , that is,

$$\mathcal{F}^+ = \{f \in \mathcal{F} : f \geq 0\} = \{f^+ : f \in \mathcal{F}\}.$$

Exercise 1.4 (Stone lattices). Show that for arbitrary $f \in \mathcal{F}$ and $c \geq 0$,

$$f \wedge c, f \vee (-c), (f - c)^+, (f + c)^- \in \mathcal{F}.$$

Find an example of a Stone lattice \mathcal{F} such that $f \vee 1 \notin \mathcal{F}$ for any $f \in \mathcal{F}$.

The following families of functions and sets play an important role in connection with integrals.

Definition 1.5 (Extension of \mathcal{F}^+ , \mathcal{F} -open sets). The *extension of \mathcal{F}^+* is the family \mathcal{F}^* of all functions $g : \Omega \rightarrow [0, \infty]$ such that

$$g = \sup_{n \geq 1} f_n$$

for some sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ .

A set $U \subset \Omega$ is called \mathcal{F} -open if

$$U = \{g > 0\}$$

for some $g \in \mathcal{F}^*$. The family of all \mathcal{F} -open sets is denoted with $\mathcal{U}(\mathcal{F})$.

These families \mathcal{F}^* and $\mathcal{U}(\mathcal{F})$ have some important properties summarized in the next three lemmas.

Lemma 1.6 (Extension of \mathcal{F}^+). For a function $g : \Omega \rightarrow [0, \infty]$, the following three properties are equivalent:

- (i) $g = \sup_{n \geq 1} f_n$ for some sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ .

- (ii) $g = \lim_{n \geq 1} f_n$ for some sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ which is pointwise increasing.
 (iii) $g = \sum_{n \geq 1} f_n$ for some sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ .

Lemma 1.7 (Properties of \mathcal{F}^*). For arbitrary scalars $c \geq 0$ and functions $g, h, g_1, g_2, g_3, \dots$ in \mathcal{F}^* , the following functions belong to \mathcal{F}^* too:

$$cg, g \wedge h, g \wedge c, \sup_{n \geq 1} g_n, \sum_{n \geq 1} g_n.$$

Lemma 1.8 (\mathcal{F} -open sets). For a set $U \subset \Omega$, the following three properties are equivalent:

- (i) $U = \{g > 0\}$ for some function $g \in \mathcal{F}^*$.
 (ii) $U = \{g > c\}$ for some function $g \in \mathcal{F}^*$ and $c \geq 0$.
 (iii) $1_U \in \mathcal{F}^*$.

Remark 1.9. Lemmas 1.7 and 1.8 imply that the family $\mathcal{U}(\mathcal{F})$ is closed under intersections and countable unions. Moreover, the σ -fields $\sigma(\mathcal{F})$ and $\sigma(\mathcal{U}(\mathcal{F}))$ coincide, where $\sigma(\mathcal{F})$ is the smallest σ -field over Ω such that all functions in \mathcal{F} are measurable, and $\sigma(\mathcal{U}(\mathcal{F}))$ is the smallest σ -field over Ω containing all sets in $\mathcal{U}(\mathcal{F})$.

To see the latter claim, note that any function in \mathcal{F}^* is a pointwise limit of a sequence in $\mathcal{F}^+ \subset \mathcal{F}$, whence $\mathcal{U}(\mathcal{F}) \subset \sigma(\mathcal{F})$, and this implies that $\sigma(\mathcal{U}(\mathcal{F})) \subset \sigma(\mathcal{F})$. On the other hand, the σ -field $\sigma(\mathcal{F})$ is generated by all sets $\{f > c\}$ and $\{f < -c\}$, $f \in \mathcal{F}$ and $c > 0$. But $\{f > c\} = \{g > 0\}$ with $g := f - f \wedge c \in \mathcal{F}^+$, so $\{f > c\} \in \mathcal{U}(\mathcal{F})$, and $\{f < -c\} = \{-f > c\}$ belongs to $\mathcal{U}(\mathcal{F})$ too. This shows that $\sigma(\mathcal{F}) \subset \sigma(\mathcal{U}(\mathcal{F}))$.

The proofs of Lemmas 1.6 and 1.7 are left to the reader as an exercise.

Proof of Lemma 1.8. It is clear that property (i) of U implies property (ii). Suppose that U has property (ii), that is, $U = \{g > c\}$ for some function $g \in \mathcal{F}^*$ and some scalar $c \geq 0$. Let $(f_n)_{n \geq 1}$ be a pointwise increasing sequence in \mathcal{F}^+ with limit g . Then

$$\tilde{f}_n := (n(f_n - c)^+) \wedge 1$$

defines a pointwise increasing sequence $(\tilde{f}_n)_{n \geq 1}$ with limit $1_U = 1_{\{g > c\}}$. Thus U has property (iii). Finally, if U has property (iii), then it has property (i) with $g := 1_U \in \mathcal{F}^*$. \square

1.2 Abstract and Usual Integrals

As before, let \mathcal{F} be a Stone lattice of functions $f : \Omega \rightarrow \mathbb{R}$ with its subcone \mathcal{F}^+ of nonnegative functions $f \in \mathcal{F}$.

Definition 1.10 (Abstract integral). A functional $J : \mathcal{F} \rightarrow \mathbb{R}$ is called an *abstract integral* (on the Stone lattice \mathcal{F}) if it has the following properties:

Linearity: For $f, g \in \mathcal{F}$ and $\lambda \in \mathbb{R}$,

$$J(\lambda f) = \lambda J(f) \quad \text{and} \quad J(f + g) = J(f) + J(g).$$

Positivity: $J(f) \geq 0$ for all $f \in \mathcal{F}^+$.

Monotone convergence: Let $(f_n)_n$ be a pointwise increasing sequence of functions $f_n \in \mathcal{F}^+$ with limit $f \in \mathcal{F}$. Then,

$$\lim_{n \rightarrow \infty} J(f_n) = J(f).$$

Linearity and positivity of an abstract integral J on \mathcal{F} imply that

$$J(f) \leq J(g) \quad \text{whenever } f, g \in \mathcal{F} \text{ with } f \leq g.$$

This follows from $g = f + h$ with $h = g - f \in \mathcal{F}^+$, so $J(g) = J(f) + J(h) \geq J(f)$.

Exercise 1.11. Let \mathcal{F} be the set of all functions (sequences) $f : \mathbb{N} \rightarrow \mathbb{R}$ such that $\lim_{\omega \rightarrow \infty} f(\omega)$ exists in \mathbb{R} . Verify that \mathcal{F} is a Stone lattice. Does $J(f) := \lim_{\omega \rightarrow \infty} f(\omega)$ define an abstract integral on \mathcal{F} ?

From now on, let $J : \mathcal{F} \rightarrow \mathbb{R}$ be an abstract integral. In view of our goal to achieve, we extend J to a functional on the extension \mathcal{F}^* of \mathcal{F}^+ . For $g \in \mathcal{F}^*$, let

$$J(g) := \sup_{f \in \mathcal{F}^+ : f \leq g} J(f).$$

Note that if $g \in \mathcal{F}^+$, then $J(f) \leq J(g)$ for all $f \in \mathcal{F}^+$ such that $f \leq g$, so the new definition of $J(g)$ yields the former value of $J(g)$. It is also obvious that for $g, h \in \mathcal{F}^*$,

$$J(g) \leq J(h) \quad \text{if } g \leq h,$$

because $g \leq h$ implies that $\{f \in \mathcal{F}^+ : f \leq g\} \subset \{f \in \mathcal{F}^+ : f \leq h\}$. The next lemma provides an alternative representation of $J(g)$ which is often convenient.

Lemma 1.12. *If $g = \lim_{n \rightarrow \infty} f_n$ with a pointwise increasing sequence $(f_n)_n$ in \mathcal{F}^+ , then $J(g) = \lim_{n \rightarrow \infty} J(f_n)$.*

Proof of Lemma 1.12. Since $(f_n)_n$ is pointwise increasing, the sequence $(J(f_n))_n$ is increasing, and $f_n \leq g$ for each n , whence $J(g) \geq \lim_{n \rightarrow \infty} J(f_n)$.

On the other hand, if f is any fixed function in \mathcal{F}^+ such that $f \leq g$, then $(f_n \wedge f)_n$ is pointwise increasing with limit f . Consequently,

$$\lim_{n \rightarrow \infty} J(f_n) \geq \lim_{n \rightarrow \infty} J(f_n \wedge f) = J(f),$$

and this shows that $J(g) \leq \lim_{n \rightarrow \infty} J(f_n)$. □

The next lemma shows that J , as a functional on \mathcal{F}^* , has various desirable properties, where we use the convention that $0 \cdot \infty := 0$.

Lemma 1.13. *For arbitrary $c \geq 0$ and functions $g, g_1, g_2, g_3, \dots \in \mathcal{F}^*$,*

$$J(cg) = cJ(g) \quad \text{and} \quad J\left(\sum_{n \geq 1} g_n\right) = \sum_{n \geq 1} J(g_n).$$

If $(g_n)_{n \geq 1}$ is pointwise increasing with limit g , then

$$J(g) = \lim_{n \rightarrow \infty} J(g_n).$$

Proof of Lemma 1.13. Let $(f_n)_{n \geq 1}$ be a pointwise increasing sequence in \mathcal{F}^+ with limit g . Then, $cf_n \in \mathcal{F}^+$ and $cf_n \uparrow cg$ as $n \rightarrow \infty$, whence

$$J(cg) = \lim_{n \rightarrow \infty} J(cf_n) = c \lim_{n \rightarrow \infty} J(f_n) = cJ(g).$$

For each $n \geq 1$, let $(f_{n,k})_{k \geq 1}$ be a pointwise increasing sequence in \mathcal{F}^+ with limit g_n . Then

$$f_N := \sum_{n=1}^N f_{n,N}$$

defines a pointwise increasing sequence $(f_N)_{N \geq 1}$ in \mathcal{F}^+ with limit $\sum_{n \geq 1} g_n$, because for any fixed integer $N_o \geq 1$ and $N \geq N_o$,

$$\sum_{n \geq 1} g_n \geq f_N \geq \sum_{n=1}^{N_o} f_{n,N} \rightarrow \sum_{n=1}^{N_o} g_n \quad \text{as } N \rightarrow \infty,$$

and $\sum_{n=1}^{N_o} g_n \uparrow \sum_{n \geq 1} g_n$ as $N_o \rightarrow \infty$. Consequently, for any fixed $N_o \geq 1$,

$$J\left(\sum_{n \geq 1} g_n\right) = \lim_{N \rightarrow \infty} J(f_N) = \lim_{N \rightarrow \infty} \sum_{n=1}^N J(f_{n,N}) \begin{cases} \leq \sum_{n \geq 1} J(g_n), \\ \geq \sum_{n=1}^{N_o} J(g_n), \end{cases}$$

and letting $N_o \rightarrow \infty$ reveals that $J(\sum_{n \geq 1} g_n) = \sum_{n \geq 1} J(g_n)$.

Suppose that $(g_n)_{n \geq 1}$ is pointwise increasing with limit g . It is tempting to write $g = \sum_{n \geq 1} \tilde{g}_n$ with $\tilde{g}_n := g_n - g_{n-1}$, $g_0 := 0$, and refer to the previous result about countable sums. But note that in general, $\tilde{g}_n \notin \mathcal{F}^*$. However,

$$f_N := \max_{n \leq N} f_{n,N}$$

defines a pointwise increasing sequence $(f_N)_{N \geq 1}$ in \mathcal{F}^+ with limit g , because for any fixed integer $N_o \geq 1$ and $N \geq N_o$,

$$g \geq g_N \geq f_N \geq \max_{n \leq N_o} f_{n,N} \uparrow g_{N_o}$$

as $N \rightarrow \infty$, and $g_{N_o} \uparrow g$ as $N_o \rightarrow \infty$. Consequently, for any fixed $N_o \geq 1$,

$$J(g) = \lim_{N \rightarrow \infty} J(f_N) \begin{cases} \leq \lim_{N \rightarrow \infty} J(g_N), \\ \geq J(g_{N_o}), \end{cases}$$

and letting $N_o \rightarrow \infty$ reveals that $J(g) = \lim_{N \rightarrow \infty} J(g_N)$. □

Now we have some essential ingredients for the proof of the following theorem.

Theorem 1.14 (Daniell–Port–Stone). *There exists a unique measure μ on the σ -field $\sigma(\mathcal{F})$ with the following properties: $\mathcal{F} \subset \mathcal{L}^1(\mu)$,*

$$J(f) = \int f d\mu \quad \text{for all } f \in \mathcal{F},$$

and for any set $B \in \sigma(\mathcal{F})$,

$$\mu(B) = \inf\{\mu(U) : U \in \mathcal{U}(\mathcal{F}), U \supset B\}$$

with $\inf(\emptyset) := \infty$.

Proof of Theorem 1.14. The proof is divided into six steps.

Step 1. For $S \subset \Omega$ let

$$\mu(S) := \inf\{J(g) : g \in \mathcal{F}^*, g \geq 1_S\}.$$

This defines an outer measure on Ω .

Proof: Since $1_\emptyset \equiv 0 \in \mathcal{F}^*$ and $J(0) = 0$, we see that $\mu(\emptyset) = 0$. To verify that μ is an outer measure, it remains to verify that for arbitrary sets $S, S_1, S_2, S_3, \dots \subset \Omega$ with $S \subset \bigcup_{n \geq 1} S_n$,

$$\mu(S) \leq \sum_{n \geq 1} \mu(S_n).$$

This is obvious if $\mu(S_n) = \infty$ for some $n \geq 1$. Otherwise, for an arbitrary fixed (small) number $\epsilon > 0$ and any index $n \geq 1$, there exists a function $g_n \in \mathcal{F}^*$ such that $g_n \geq 1_{S_n}$ and $J(g_n) \leq \mu(S_n) + 2^{-n}\epsilon$. But then the definition of μ and Lemma 1.13 imply that

$$\mu(S) \leq J\left(\sup_{n \geq 1} g_n\right) \leq J\left(\sum_{n \geq 1} g_n\right) = \sum_{n \geq 1} J(g_n) \leq \sum_{n \geq 1} \mu(S_n) + \epsilon.$$

As $\epsilon \downarrow 0$, we obtain the desired inequality.

Step 2. $\mu(U) = J(1_U)$ for any $U \in \mathcal{U}(\mathcal{F})$.

Proof: By Lemma 1.8, $1_U \in \mathcal{F}^*$, and $J(g) \geq J(1_U)$ for any $g \in \mathcal{F}^*$ such that $g \geq 1_U$.

Step 3. Every set $U \in \mathcal{U}(\mathcal{F})$ is μ -measurable, that is,

$$\mu(S) = \mu(S \cap U) + \mu(S \setminus U).$$

In particular, by Carathéodory's theory of outer measures, the restriction of μ to $\sigma(\mathcal{F})$ defines a measure on that σ -field.

Proof: Since μ is an outer measure, it suffices to show that $\mu(S) \geq \mu(S \cap U) + \mu(S \setminus U)$ in case of $\mu(S) < \infty$. Let $g \in \mathcal{F}^*$ with $g \geq 1_S$. Let $(g_n)_{n \geq 1}$ and $(f_n)_{n \geq 1}$ be pointwise increasing sequences in \mathcal{F}^+ with limits g and $g \wedge 1_U$, respectively. Without loss of generality let $g_n \geq f_n$ for all $n \geq 1$. On the one hand,

$$\mu(S \cap U) \leq J(g \wedge 1_U) = \lim_{n \rightarrow \infty} J(f_n).$$

On the other hand, for any fixed $n_o \geq 1$,

$$1_{S \setminus U} = 1_S(1 - 1_U) \leq g(1 - 1_U) \leq g - g \wedge 1_U \leq g - f_{n_o}.$$

But $g - f_{n_o}$ is the limit of the pointwise increasing sequence $(g_n - f_{n_o})_{n \geq n_o}$ in \mathcal{F}^+ , whence $g - f_{n_o} \in \mathcal{F}^*$ and

$$\begin{aligned} \mu(S \setminus U) &\leq J(g - f_{n_o}) = \lim_{n \rightarrow \infty} J(g_n - f_{n_o}) \\ &= \lim_{n \rightarrow \infty} J(g_n) - J(f_{n_o}) \\ &= J(g) - J(f_{n_o}). \end{aligned}$$

Consequently, $\mu(S \cap U) + \mu(S \setminus U) \leq J(g) + \lim_{n \rightarrow \infty} J(f_n) - J(f_{n_o})$, and as $n_o \rightarrow \infty$, we obtain the inequality

$$\mu(S \cap U) + \mu(S \setminus U) \leq J(g)$$

for any $g \in \mathcal{F}^*$ such that $g \geq 1_S$. This implies the desired inequality $\mu(S \cap U) + \mu(S \setminus U) \leq \mu(S)$.

Step 4. $\mathcal{F} \subset \mathcal{L}^1(\mu)$ and $J(f) = \int f d\mu$ for all $f \in \mathcal{F}$.

Proof: It suffices to show that $J(f) = \int f d\mu$ for any $f \in \mathcal{F}^+$. For $n \geq 1$ let

$$f_n := 2^{-n} \sum_{k \geq 1} 1_{\{f > 2^{-n}k\}}.$$

For any $\omega \in \{f > 0\}$, $f_n(\omega)$ is the largest point on the grid $\{2^{-n}z : z \in \mathbb{N}_0\}$ which is strictly smaller than $f(\omega)$. Since $\{2^{-n}z : z \in \mathbb{N}_0\} \subset \{2^{-(n+1)}z : z \in \mathbb{N}_0\}$ and $f_n = 0$ on $\{f = 0\}$, the function sequence $(f_n)_{n \geq 1}$ is pointwise increasing with limit f . For any constant $c \geq 0$, the set $\{f > c\}$ is \mathcal{F} -open, and

$$J(1_{\{f > c\}}) = \mu(\{f > c\}) = \int 1_{\{f > c\}} d\mu.$$

Consequently, $f_n \in \mathcal{F}^*$, so by Lemma 1.13 and monotone convergence for usual integrals,

$$J(f) = \lim_{n \rightarrow \infty} J(f_n) \quad \text{and} \quad \int f d\mu = \lim_{n \rightarrow \infty} \int f_n d\mu.$$

But for any fixed n , Lemma 1.13 and standard properties of usual integrals reveal that

$$J(f_n) = 2^{-n} \sum_{k \geq 1} J(1_{\{f > 2^{-n}k\}}) = 2^{-n} \sum_{k \geq 1} \int 1_{\{f > 2^{-n}k\}} d\mu = \int f_n d\mu.$$

This shows that $J(f) = \int f d\mu$.

Step 5. For any set $S \subset \Omega$,

$$\mu(S) = \tilde{\mu}(S) := \inf \{ \mu(U) : U \in \mathcal{U}(\mathcal{F}), U \supset S \}.$$

Proof: Since $U \supset S$ is equivalent to $1_U \geq 1_S$, and since $1_U \in \mathcal{F}^*$ for any $U \in \mathcal{U}(\mathcal{F})$, we obtain the inequality $\mu(S) \leq \tilde{\mu}(S)$.

On the other hand, for any fixed number $c \in (0, 1)$ and arbitrary functions $g \in \mathcal{F}^*$ with $g \geq S$,

$$c^{-1}g \geq 1_{\{g > c\}} \geq 1_S.$$

Since $\{g > c\} \in \mathcal{U}(\mathcal{F})$, this implies that

$$J(g) \geq c\mu(\{g > c\}) \geq c\tilde{\mu}(S).$$

Letting $c \uparrow 1$ reveals that $J(g) \geq \tilde{\mu}(S)$, whence $\mu(S) \geq \tilde{\mu}(S)$.

Step 6. A measure μ on $\sigma(\mathcal{F})$ with the stated properties is unique.

Proof: The measure μ is uniquely determined by its values $\mu(U)$, $U \in \mathcal{U}(\mathcal{F})$. But for each $U \in \mathcal{U}(\mathcal{F})$ there exists a pointwise increasing sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ with limit 1_U , so

$$\mu(U) = \int 1_U d\mu = \lim_{n \rightarrow \infty} \int f_n d\mu = \lim_{n \rightarrow \infty} J(f_n).$$

Thus, the measure μ is uniquely determined by the given functional $J : \mathcal{F} \rightarrow \mathbb{R}$. \square

1.3 Representations of Dual Spaces

Suppose for the moment that \mathcal{F} is an arbitrary real vector space, equipped with a seminorm $\|\cdot\|$. The dual space of $(\mathcal{F}, \|\cdot\|)$ is the space of linear functionals $L : \mathcal{F} \rightarrow \mathbb{R}$ which are continuous with respect to the seminorm $\|\cdot\|$. Continuity of L is equivalent to

$$\sup_{f \in \mathcal{F} : \|f\| \leq 1} |L(f)| < \infty.$$

In Functional Analysis, an important question is how the dual space can be represented explicitly. Many of the known results are called a Riesz representation, honoring the hungarian mathematician Frigyes Riesz (1880-1956) who (co-)derived many of these results.

Specifically, let \mathcal{F} be a Stone lattice of functions on a set Ω . Suppose in addition that there is a seminorm $\|\cdot\|$ on \mathcal{F} satisfying the following two properties:

$$\|f\| \leq \|g\| \quad \text{whenever } f, g \in \mathcal{F}^+ \text{ with } f \leq g.$$

For any pointwise decreasing sequence $(f_n)_{n \geq 1}$ in \mathcal{F}^+ with limit 0,

$$\lim_{n \rightarrow \infty} \|f_n\| = 0.$$

The next theorem shows that any continuous linear functional on $(\mathcal{F}, \|\cdot\|)$ can be represented with certain integrals.

Theorem 1.15. *For any functional L in the dual space of $(\mathcal{F}, \|\cdot\|)$ there exist measures μ^+, μ^- on $\sigma(\mathcal{F})$ such that $\mathcal{F} \subset \mathcal{L}^1(\mu^+ + \mu^-)$ and*

$$L(f) = \int f d\mu^+ - \int f d\mu^- \quad \text{for all } f \in \mathcal{F}.$$

Here is a first special case of this theorem:

Theorem 1.16 (Riesz–Markov–Kakutani). *Let (Ω, d) be a compact metric space, and let $\mathcal{C}(\Omega)$ be the family of continuous functions $f : \Omega \rightarrow \mathbb{R}$ with respect to d , equipped with the supremum norm $\|\cdot\|_\infty$, that is, $\|f\|_\infty = \max_{\omega \in \Omega} |f(\omega)|$. Let $L : \mathcal{C}(\Omega) \rightarrow \mathbb{R}$ be a linear functional which is continuous with respect to $\|\cdot\|_\infty$. Then there exist finite measures μ^+, μ^- on $\text{Borel}(\Omega)$ such that*

$$L(f) = \int f d\mu^+ - \int f d\mu^- \quad \text{for all } f \in \mathcal{C}(\Omega).$$

Once we have learned more about signed measures, the conclusion of Theorem 1.16 can be reformulated as follows: There exists a probability measure P on $\text{Borel}(\Omega)$ and a bounded, measurable function $h: \Omega \rightarrow \mathbb{R}$ such that

$$L(f) = \int fh \, dP \quad \text{for all } f \in \mathcal{C}(\Omega).$$

The proof of Theorem 1.15 uses the following two auxiliary results.

Lemma 1.17 (From additive, nonnegative to linear functionals). *Suppose that $J: \mathcal{F}^+ \rightarrow [0, \infty)$ is a functional which is additive in the sense that $J(f + g) = J(f) + J(g)$ for all $f, g \in \mathcal{F}^+$. Then,*

$$L(f) := J(f^+) - J(f^-)$$

defines a linear functional $L: \mathcal{F} \rightarrow \mathbb{R}$ such that $L \equiv J$ on \mathcal{F}^+ .

Lemma 1.18 (From linear to additive, nonnegative functionals). *Suppose that $L: \mathcal{F} \rightarrow \mathbb{R}$ is a linear functional. For $f \in \mathcal{F}^+$ let*

$$\begin{aligned} J(f) &:= \sup\{L(h): h \in \mathcal{F}^+, h \leq f\}, \\ K(f) &:= \sup\{-L(h): h \in \mathcal{F}^+, h \leq f\}. \end{aligned}$$

These two functionals $J, K: \mathcal{F}^+ \rightarrow [0, \infty]$ are additive in the sense that $J(f + g) = J(f) + J(g)$ and $K(f + g) = K(f) + K(g)$ for all $f, g \in \mathcal{F}^+$.

They satisfy the equations $J = K + L$ and $K = J - L$. In particular, J is real-valued if and only if K is real-valued.

If J and K are real-valued, then for arbitrary $f \in \mathcal{F}$,

$$L(f) = L_J(f) - L_K(f),$$

where $L_J(f) := J(f^+) - J(f^-)$ and $L_K(f) := K(f^+) - K(f^-)$.

The proof of Lemma 1.17 is left to the reader as an exercise.

Proof of Lemma 1.18. Since $J(f) \geq L(0) = 0$, the functional J is nonnegative. Let $f_1, f_2 \in \mathcal{F}^+$. For arbitrary $h_1, h_2 \in \mathcal{F}^+$ with $h_1 \leq f_1$ and $h_2 \leq f_2$,

$$J(f_1 + f_2) \geq L(h_1 + h_2) = L(h_1) + L(h_2),$$

and letting $L(h_j) \uparrow J(f_j)$ for $j = 1, 2$ reveals that

$$J(f_1 + f_2) \geq J(f_1) + J(f_2).$$

On the other hand, if $h \in \mathcal{F}^+$ with $h \leq f_1 + f_2$, then $h = h_1 + h_2$ with the functions

$$h_1 := h \wedge f_1 \leq f_1 \quad \text{and} \quad h_2 := (h - f_1)^+ \leq f_2$$

in \mathcal{F}^+ . Hence,

$$L(h) = L(h_1) + L(h_2) \leq J(f_1) + J(f_2),$$

and letting $L(h) \uparrow J(f_1 + f_2)$ reveals that

$$J(f_1 + f_2) \leq J(f_1) + J(f_2).$$

Replacing L with $-L$ reveals that K is also a nonnegative, additive functional on \mathcal{F}^+ .

For $f \in \mathcal{F}^+$,

$$\begin{aligned} J(f) &= \sup\{L(h) : h \in \mathcal{F}^+, h \leq f\} \\ &= \sup\{L(f) - L(f - h) : h \in \mathcal{F}^+, h \leq f\} \\ &= L(f) + \sup\{-L(f - h) : h \in \mathcal{F}^+, h \leq f\} \\ &= L(f) + \sup\{-L(h) : h \in \mathcal{F}^+, h \leq f\} \\ &= L(f) + K(f), \end{aligned}$$

where the second to last step uses the fact that the sets $\{f - h : h \in \mathcal{F}^+, h \leq f\}$ and $\{h \in \mathcal{F}^+ : h \leq f\}$ coincide. Thus, $J = K + L$ on \mathcal{F}^+ , and this is equivalent to $K = J - L$.

Now suppose that J and K are real-valued on \mathcal{F}^+ . Then we know from Lemma 1.17 that $L_J(f) := J(f^+) - J(f^-)$ and $L_K(f) := K(f^+) - K(f^-)$ define linear functionals on \mathcal{F} . Since $L = J - K$ on \mathcal{F}^+ ,

$$L(f) = L(f^+) - L(f^-) = J(f^+) - K(f^+) - J(f^-) + K(f^-) = L_J(f) - L_K(f)$$

for any function $f \in \mathcal{F}$. □

Proof of Theorem 1.15. Continuity of L means that $|L(f)| \leq C\|f\|$ for all $f \in \mathcal{F}$ and some real constant $C = C(L) \geq 0$. The additional properties of the seminorm $\|\cdot\|$ imply that for $f \in \mathcal{F}^+$,

$$\begin{aligned} J(f) &:= \sup\{L(h) : h \in \mathcal{F}^+, h \leq f\} \begin{cases} \geq \max\{0, L(f)\}, \\ \leq C\|f\|, \end{cases} \\ K(f) &:= \sup\{-L(h) : h \in \mathcal{F}^+, h \leq f\} \begin{cases} \geq \max\{0, -L(f)\}, \\ \leq C\|f\|. \end{cases} \end{aligned}$$

Consequently, J and K are additive functionals \mathcal{F}^+ with values in $[0, \infty)$, and L_J, L_K are well-defined linear functionals on \mathcal{F} with $L = L_J - L_K$, see Lemmas 1.18 and 1.17. The sequential continuity property of $\|\cdot\|$ implies that L_J, L_K are abstract integrals on \mathcal{F} . Indeed, if $(f_n)_{n \geq 1}$ is a pointwise increasing sequence in \mathcal{F}^+ with limit $f \in \mathcal{F}^+$, then for $M = J, K$,

$$L_M(f_n) = L_M(f) - L_M(f - f_n) = L_M(f) - M(f - f_n) \begin{cases} \geq L_M(f) - C\|f - f_n\|, \\ \leq L_M(f). \end{cases}$$

As $n \rightarrow \infty$, $f - f_n \downarrow 0$, whence $\|f - f_n\| \rightarrow 0$ and, consequently, $L_M(f_n) \rightarrow L_M(f)$. According to Theorem 1.14, there exist measures μ^+, μ^- on $\mathcal{U}(\mathcal{F})$ such that for arbitrary $f \in \mathcal{F}$, $L_J(f) = \int f d\mu^+$ and $L_K(f) = \int f d\mu^-$, which implies the asserted representation of L . □

The proof of Theorem 1.16 uses a well-known result from analysis about pointwise and uniform convergence.

Lemma 1.19 (Dini). *Let (Ω, d) be a compact metric space, and let $(f_n)_{n \geq 1}$ be a pointwise decreasing sequence of continuous functions with limit 0. Then $\|f_n - f\|_\infty \rightarrow 0$ as $n \rightarrow \infty$.*

Proof of Lemma 1.19. For any $\epsilon > 0$, the sets $\{f_n < \epsilon\}$ are open subsets of Ω with $\{f_1 < \epsilon\} \subset \{f_2 < \epsilon\} \subset \{f_3 < \epsilon\} \subset \dots$, and $\Omega = \bigcup_{n \geq 1} \{f_n < \epsilon\}$. But compactness of Ω implies that $\Omega = \bigcup_{n \leq n(\epsilon)} \{f_n < \epsilon\} = \{f_{n(\epsilon)} < \epsilon\}$ for some integer $n(\epsilon) \geq 1$, and this implies that $\|f_n\|_\infty < \epsilon$ for all $n \geq n(\epsilon)$. \square

Proof of Theorem 1.16. In view of Theorem 1.15, note first that the supremum norm $\|\cdot\|_\infty$ has the additional required properties. It is obvious that $\|f\|_\infty \leq \|g\|_\infty$ for functions $0 \leq f \leq g$, and the second property about sequences follows from Lemma 1.19. Thus we can apply Theorem 1.15 and conclude that there exist measures μ^+ and μ^- on $\sigma(\mathcal{C}(\Omega))$ such that $\mathcal{F} \subset \mathcal{L}^1(\mu^+ + \mu^-)$ and $L(f) = \int f d\mu^+ - \int f d\mu^-$. Since $1 \in \mathcal{F}$, these measures μ^\pm are finite. Moreover, $\sigma(\mathcal{C}(\Omega))$ coincides with $\text{Borel}(\Omega)$. On the one hand, $\{f > r\}$ is an open subset of Ω for all $f \in \mathcal{C}(\Omega)$ and $r \in \mathbb{R}$, whence $\sigma(\mathcal{C}(\Omega)) \subset \text{Borel}(\Omega)$. On the other hand, if U is a nonvoid proper subset of Ω , then $f(x) := \inf\{d(x, y) : y \in \Omega \setminus U\}$ defines a continuous function such that $U = \{f > 0\}$. Hence, $\text{Borel}(\Omega) \subset \sigma(\mathcal{C}(\Omega))$. \square

Chapter 2

Signed Measures

2.1 The Hahn–Jordan Decomposition

For simplicity we restrict our attention to finite signed measures. Throughout this section let (Ω, \mathcal{A}) be a measurable space. That means, Ω is a nonvoid set equipped with a σ -field \mathcal{A} over Ω , the family of measurable subsets of Ω .

Definition 2.1 (Finite signed measure). A function $\nu : \mathcal{A} \rightarrow \mathbb{R}$ is called a *finite signed measure* on (Ω, \mathcal{A}) , if

(SM.1) $\nu(\emptyset) = 0$ and

(SM.2) ν is σ -additive, that is, for arbitrary disjoint sets A_1, A_2, A_3, \dots in \mathcal{A} ,

$$\nu\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} \nu(A_n).$$

Note that the only difference to a finite measure is that $\nu(A)$ may be negative for some sets $A \in \mathcal{A}$.

Example 2.2. A standard example of a finite signed measure is $\nu := Q - P$ with finite measures P, Q on (Ω, \mathcal{A}) . Properties (SM.1-2) follow immediately from analogous properties of measures.

Example 2.3. Another example is given by

$$\nu(A) := \int_A f d\mu = \int 1_A f d\mu$$

with a measure μ on (Ω, \mathcal{A}) and a function $f \in \mathcal{L}^1(\mu)$. That means, $f : \Omega \rightarrow \bar{\mathbb{R}}$ is \mathcal{A} -measurable with $\int |f| d\mu < \infty$. Here one can verify properties (SM.1-2) by means of linearity of integrals and dominated convergence.

Later on it will be shown that any finite signed measure may be represented as in Examples 2.2 and 2.3.

Remark 2.4 (Additivity). A finite signed measure ν is additive in the sense that

$$\nu\left(\bigcup_{n=1}^N A_n\right) = \sum_{n=1}^N \nu(A_n)$$

for arbitrary $N \in \mathbb{N}$ and disjoint sets $A_1, \dots, A_N \in \mathcal{A}$. This follows from (SM.1-2) by setting $A_n = \emptyset$ for $n > N$.

Exercise 2.5 (Continuity properties of signed measures). Let ν be a finite signed measure on (Ω, \mathcal{A}) . Show that for arbitrary sets $B_1 \subset B_2 \subset B_3 \subset \dots$ in \mathcal{A} ,

$$\nu\left(\bigcup_{n=1}^{\infty} B_n\right) = \lim_{n \rightarrow \infty} \nu(B_n).$$

Show that for arbitrary sets $C_1 \supset C_2 \supset C_3 \supset \dots$ in \mathcal{A} ,

$$\nu\left(\bigcap_{n=1}^{\infty} C_n\right) = \lim_{n \rightarrow \infty} \nu(C_n).$$

The next definition introduces two important concepts for the subsequent results.

Definition 2.6 (Positive and negative sets). Let ν be a finite signed measure on (Ω, \mathcal{A}) . A set $A_* \subset \Omega$ is called ν -positive if $A_* \in \mathcal{A}$ and

$$\nu(A) \geq 0 \quad \text{for all } A \in \mathcal{A} \text{ with } A \subset A_*.$$

A set $A_* \subset \Omega$ is called ν -negative, if $A_* \in \mathcal{A}$ and

$$\nu(A) \leq 0 \quad \text{for all } A \in \mathcal{A} \text{ with } A \subset A_*.$$

Since $-\nu$ is a finite signed measure too, a set $A_* \subset \Omega$ is ν -positive or ν -negative if and only if it is $(-\nu)$ -negative or $(-\nu)$ -positive, respectively.

Here is a key result for the main theorems in this and the next section.

Proposition 2.7 (Existence of nontrivial positive sets). *Let ν be a finite signed measure on (Ω, \mathcal{A}) , and let $A_0 \in \mathcal{A}$ with $\nu(A_0) > 0$. Then there exists a ν -positive set $A_* \subset A_0$ with $\nu(A_*) \geq \nu(A_0)$.*

Proof of Proposition 2.7. We define

$$\delta_0 := \sup\{-\nu(B) : B \in \mathcal{A}, B \subset A_0\} \geq 0.$$

Then we write $A_0 = A_1 \cup B_1$ with disjoint measurable sets A_1 and B_1 such that

$$-\nu(B_1) \geq \min\{\delta_0/2, 1\}.$$

This procedure can be iterated. After k steps we have measurable sets $A_0 \supset A_1 \supset \dots \supset A_k$, and we consider the number

$$\delta_k := \sup\{-\nu(B) : B \in \mathcal{A}, B \subset A_k\} \geq 0.$$

Then we write $A_k = A_{k+1} \cup B_{k+1}$ with disjoint measurable sets A_{k+1} and B_{k+1} such that

$$-\nu(B_{k+1}) \geq \min\{\delta_k/2, 1\}.$$

This construction yields a non-increasing sequence $(A_k)_{k \geq 0}$ and the disjoint sets $B_k = A_{k-1} \setminus A_k$, $k \geq 1$. We may write

$$A_0 = A_* \cup B_*$$

with the disjoint sets

$$A_* = \bigcap_{k \geq 0} A_k \quad \text{and} \quad B_* = \bigcup_{k \geq 1} B_k.$$

Since $\nu(B_*) = \sum_{k=1}^{\infty} \nu(B_k)$ with nonpositive summands $\nu(B_k)$, we obtain the inequality

$$\nu(A_*) = \nu(A_0) - \nu(B_*) \geq \nu(A_0).$$

Moreover, the sequence $(\nu(B_k))_{k \geq 1}$ converges to 0, so the inequalities $0 \leq \min\{\delta_k/2, 1\} \leq -\nu(B_{k+1})$ imply that $\lim_{k \rightarrow \infty} \delta_k = 0$. Consequently, for any measurable set $A \subset A_*$,

$$-\nu(A) \leq \inf_{k \geq 0} \sup\{-\nu(A') : A' \in \mathcal{A}, A' \subset A_k\} = \inf_{k \geq 0} \delta_k = 0,$$

whence $\nu(A) \geq 0$. This shows that A_* is a ν -positive set. \square

Exercise 2.8 (Unions of positive sets). Let ν be a finite signed measure on (Ω, \mathcal{A}) . Let $(A_n)_{n \geq 1}$ be a sequence of ν -positive sets. Show that $A_* := \bigcup_{n \geq 1} A_n$ is also ν -positive and satisfies

$$\nu(A_*) \geq \sup_{n \geq 1} \nu(A_n).$$

With Proposition 2.7 and Exercise 2.8 one can prove the following representation of signed measures:

Theorem 2.9 (Hahn–Jordan decomposition). *Let ν be a finite signed measure on (Ω, \mathcal{A}) . Then $\Omega = \Omega_+ \cup \Omega_-$ with disjoint sets Ω_+, Ω_- such that Ω_+ is ν -positive and Ω_- is ν -negative. In other words,*

$$\nu^+(A) := \nu(A \cap \Omega_+) \quad \text{and} \quad \nu^-(A) := -\nu(A \cap \Omega_-)$$

defines two finite measures ν^+, ν^- on (Ω, \mathcal{A}) such that $\nu = \nu^+ - \nu^-$, and these two measures have disjoint support in the sense that $\nu^+(\Omega_-) = \nu^-(\Omega_+) = 0$.

Example 2.10. Before proving Theorem 2.9, we revisit Example 2.3, i.e. $\nu(A) = \int_A f d\mu$ for some measure μ on (Ω, \mathcal{A}) and a function $f \in \mathcal{L}^1(\mu)$. Here a Hahn–Jordan decomposition is given by

$$\Omega_+ := \{f \geq 0\} \quad \text{and} \quad \Omega_- := \{f < 0\}.$$

The corresponding measures ν^\pm are given by

$$\nu^\pm(A) = \int_A f^\pm d\mu$$

with $f^\pm(\omega) = \max\{\pm f(\omega), 0\}$.

Remark 2.11. Theorem 2.9 implies that any finite signed measure ν may be represented as in Example 2.3. Indeed, let $\nu, \Omega_+, \Omega_-, \nu^+$ and ν^- be as in Theorem 2.9. Then

$$\nu(A) = \int_A f d\mu \quad \text{for all } A \in \mathcal{A},$$

with the finite measure $\mu := \nu^+ + \nu^-$ and the measurable function $f := 1_{\Omega_+} - 1_{\Omega_-}$.

Proof of Theorem 2.9. Let $(A_n)_{n \geq 1}$ be a sequence of measurable sets such that

$$\lim_{n \rightarrow \infty} \nu(A_n) = C_+ := \sup\{\nu(A) : A \in \mathcal{A}\}.$$

According to Proposition 2.7, we may assume without loss of generality that all sets A_n are ν -positive. But then,

$$\Omega_+ := \bigcup_{n \geq 1} A_n$$

is a ν -positive set such that $\nu(\Omega_+) = C_+$, see Exercise 2.8. In particular, $C_+ < \infty$. Moreover, $\Omega_- := \Omega \setminus \Omega_+$ is a ν -negative set. Indeed, if $\nu(A_0) > 0$ for some measurable set $A_0 \subset \Omega_-$, then $\nu(\Omega_+ \cup A_0) = \nu(\Omega_+) + \nu(A_0) > C_+$, a contradiction to the definition of C_+ . \square

Exercise 2.12. Let $\Omega = \Omega_+ \cup \Omega_-$ as in Theorem 2.9. Show that this partition of Ω is essentially unique in the following sense: If $\Omega = B_+ \cup B_-$ with a ν -positive set B_+ and a ν -negative set B_- , then

$$\nu(A) = 0 \quad \text{for any measurable } A \subset (\Omega_+ \Delta B_+) \cup (\Omega_- \Delta B_-).$$

Remark 2.13. Theorem 2.9 implies, that any finite signed measure ν may be represented as the difference of two finite measures. The next exercise shows that the particular measures ν^\pm are minimal in a certain sense.

Exercise 2.14. Let ν be a finite signed measure on (Ω, \mathcal{A}) with Hahn–Jordan decomposition $\nu = \nu^+ - \nu^-$.

(a) Show that $-\nu^-(A) \leq \nu(A) \leq \nu^+(A)$ for arbitrary $A \in \mathcal{A}$.

(b) Let $\nu = \mu^+ - \mu^-$ with finite measures μ^+, μ^- on (Ω, \mathcal{A}) . Show that there exists a finite measure μ_o on (Ω, \mathcal{A}) such that $\mu^+ = \nu^+ + \mu_o$ and $\mu^- = \nu^- + \mu_o$. Deduce from the latter fact that ν^+ and ν^- are “minimal” in the sense that

$$\mu^+ + \mu^- \geq \nu^+ + \nu^-$$

with equality if and only if $\mu^+ = \nu^+$ and $\mu^- = \nu^-$.

Exercise 2.15. Let P and Q be finite measures on \mathbb{R} with densities f and g , respectively, with respect to Lebesgue measure. Determine a Hahn–Jordan decomposition of $\nu := Q - P$ in terms of f and g .

Illustrate your solution graphically in case of $P = \mathcal{N}(0, 1)$ and $Q = \mathcal{N}(0, 2^2)$.

2.2 Radon–Nikodym Derivatives

In this section we consider measures P and Q on (Ω, \mathcal{A}) and investigate under which conditions Q has a density f with respect to P . That means, $f : \Omega \rightarrow [0, \infty)$ is an \mathcal{A} -measurable function such that

$$(2.1) \quad Q(A) = \int_A f dP \quad \text{for all } A \in \mathcal{A}.$$

Such a function f is also called a *Radon–Nikodym derivative of Q with respect to P* , and (2.1) is sometimes abbreviated as

$$f = \frac{dQ}{dP}.$$

If a density of Q with respect to P exists, arbitrary integrals with respect to Q may be rewritten as integrals with respect to P .

Lemma 2.16. *Suppose that Q has a density f with respect to P in the sense of (2.1). Then for arbitrary \mathcal{A} -measurable functions $h : \Omega \rightarrow \bar{\mathbb{R}}$,*

$$(2.2) \quad \int h dQ = \int hf dP$$

whenever one of the two integrals is well-defined in $\bar{\mathbb{R}}$.

Proof of Lemma 2.16. We follow a standard route in measure theory, also known as “measure-theoretic induction”. Assumption (2.1) is equivalent to (2.2) for arbitrary indicator functions $h = 1_A$, $A \in \mathcal{A}$. By linearity of integrals, (2.2) holds true for “simple functions” h , that means, functions $h = \sum_{j=1}^m \lambda_j 1_{A_j}$ with $m \in \mathbb{N}$, constants $\lambda_j \geq 0$ and sets $A_j \in \mathcal{A}$. For an arbitrary measurable function $h : \Omega \rightarrow [0, \infty]$, there exists a sequence $(h_n)_n$ of simple functions such that $(h_n)_n \uparrow h$ pointwise. A standard construction is given by

$$h_n(\omega) := 2^{-n} \sum_{j=1}^{n2^n} 1_{[h(\omega) \geq 2^{-n}j]}.$$

But then $(h_n f)_n \uparrow hf$, so by monotone convergence,

$$\int h dQ = \lim_{n \rightarrow \infty} \int h_n dQ = \lim_{n \rightarrow \infty} \int h_n f dP = \int hf dP.$$

Finally, any measurable function $h : \Omega \rightarrow \bar{\mathbb{R}}$ may be written as $h = h^+ - h^-$ with $h^\pm = (\pm h) \vee 0$. Then $(hf)^\pm = h^\pm f$, because $f \geq 0$, and $\int h^\pm dQ = \int (hf)^\pm dP$. Hence

$$\int h dQ = \int h^+ dQ - \int h^- dQ = \int (hf)^+ dP - \int (hf)^- dP = \int hf dP,$$

whenever these differences are well-defined in $\bar{\mathbb{R}}$. □

Corollary 2.17. *Let P , Q and R be measures on (Ω, \mathcal{A}) such that densities $f = dQ/dP$ and $g = dR/dQ$ exist. Then $fg = dR/dP$. Furthermore, if $f > 0$, then $f^{-1} = dP/dQ$.*

Proof of Corollary 2.17. For any set $A \in \mathcal{A}$,

$$\int_A fg dP = \int 1_A g f dP = \int 1_A g dQ = \int_A g dQ = R(A),$$

where the second step follows from Lemma 2.16 applied to $h = 1_A g$.

In case of $f > 0$, for any $A \in \mathcal{A}$,

$$P(A) = \int 1_A dP = \int 1_A f^{-1} f dP = \int 1_A f^{-1} dQ = \int_A f^{-1} dQ$$

by Lemma 2.16 applied to $h = 1_A f^{-1}$. □

Exercise 2.18. Let P be a measure on (Ω, \mathcal{A}) , and for a fixed $\Omega_o \in \mathcal{A}$ let $P_o(A) := P(A \cap \Omega_o)$. Show that P_o defines a measure on (Ω, \mathcal{A}) with $P_o(\Omega) = P(\Omega_o)$ and $P_o(\Omega \setminus \Omega_o) = 0$. Determine a version of dP_o/dP . Then deduce that for arbitrary measurable functions $h : \Omega \rightarrow \bar{\mathbb{R}}$,

$$\int h dP_o = \int_{\Omega_o} h dP$$

whenever one (and thus both) of these integrals is well-defined in $\bar{\mathbb{R}}$.

Exercise 2.19 (Transformations). Let (Ω, \mathcal{A}) and $(\tilde{\Omega}, \tilde{\mathcal{A}})$ be measurable spaces, and let $\tau : \Omega \rightarrow \tilde{\Omega}$ be a bijective mapping such that τ and τ^{-1} are measurable.

(a) Let M be a measure on (Ω, \mathcal{A}) , and let $\tilde{M} := M \circ \tau^{-1}$, a measure on $(\tilde{\Omega}, \tilde{\mathcal{A}})$. Show that

$$\int \tilde{h} d\tilde{M} = \int \tilde{h} \circ \tau dM$$

for arbitrary measurable functions $\tilde{h} : \tilde{\Omega} \rightarrow \bar{\mathbb{R}}$ such that one of these integrals is well-defined in $\bar{\mathbb{R}}$.

(b) Let P and Q be measures on (Ω, \mathcal{A}) such that a density $f = dQ/dP$ exists. Show that $\tilde{Q} := Q \circ \tau^{-1}$ has density $\tilde{f} := f \circ \tau^{-1}$ with respect to $\tilde{P} := P \circ \tau^{-1}$.

2.2.1 Finite measures

Now we consider the case of finite measures P and Q in more detail. The next lemma implies that a density of Q with respect to P , if it exists, is P -almost everywhere unique.

Lemma 2.20 (Uniqueness of densities). *Let P be a finite measure on (Ω, \mathcal{A}) , and for $j = 1, 2$ let $Q_j : \mathcal{A} \rightarrow \mathbb{R}$ be given by $Q_j(A) := \int_A f_j dP$ with a function $f_j \in \mathcal{L}^1(P)$. If $Q_1 \leq Q_2$ on \mathcal{A} , then*

$$P(\{f_1 > f_2\}) = 0.$$

Exercise 2.21. Prove Lemma 2.20.

In what follows, we shall *construct* a density of Q with respect to P . The main idea for the construction results from an elementary consideration in the next exercise.

Exercise 2.22 (Superlevel sets). Let P and Q be finite measures on (Ω, \mathcal{A}) such that Q has a density $f \in \mathcal{L}^1(P)$ with respect to P . Show that for any $\lambda > 0$, a set $A \in \mathcal{A}$ with $\{f > \lambda\} \subset A \subset \{f \geq \lambda\}$ maximizes

$$\mathcal{A} \ni A \mapsto Q(A) - \lambda P(A).$$

More generally, show that a set $A \in \mathcal{A}$ maximizes $Q - \lambda P$ if and only if

$$P(A \cap \{f < \lambda\}) = 0 = P(\{f > \lambda\} \setminus A).$$

Now we state and prove the first main result of this section.

Theorem 2.23. *Let P and Q be finite measures on (Ω, \mathcal{A}) . There exist a set $B_* \in \mathcal{A}$ with $P(B_*) = 0$ and a nonnegative function $f \in \mathcal{L}^1(P)$ such that*

$$Q(A) = Q(A \cap B_*) + \int_A f dP \quad \text{for arbitrary } A \in \mathcal{A}.$$

Proof of Theorem 2.23. Motivated by Exercise 2.22, we consider for any $\lambda > 0$ the finite signed measure $Q - \lambda P$. According to Theorem 2.9, there exists a set $A_\lambda \in \mathcal{A}$ such that A_λ is $(Q - \lambda P)$ -positive and $\Omega \setminus A_\lambda$ is $(Q - \lambda P)$ -negative. Now we could try to combine all these sets and define a function f via $\{f > \lambda\} = A_\lambda$. But this is easier said than done. The problem is that we are dealing with uncountably many sets A_λ .

The precise construction starts from a countable dense subset Λ of $(0, \infty)$, for instance, $\Lambda = \mathbb{Q} \cap (0, \infty)$. For any $\lambda \in \Lambda$ let $A_\lambda \in \mathcal{A}$ be chosen as above. Now we “clean up” these sets as follows: For $t \geq 0$ let

$$B_t := \bigcup_{\lambda \in \Lambda: \lambda > t} A_\lambda.$$

Note that $\{A_\lambda : \lambda \in \Lambda, \lambda > t\}$ is a countable family of $(Q - tP)$ -positive sets, because $Q - \lambda P \leq Q - tP$ for $\lambda > t$. Hence, B_t is measurable and $(Q - tP)$ -positive, too, see Exercise 2.8. Moreover, for any measurable subset A of $\Omega \setminus B_t = \bigcap_{\lambda \in \Lambda: \lambda > t} \Omega \setminus A_\lambda$,

$$(Q - tP)(A) = \sup_{\lambda \in \Lambda: \lambda > t} (Q - \lambda P)(A) \leq 0,$$

because $\Omega \setminus A_\lambda$ is $(Q - \lambda P)$ -negative. Hence $\Omega \setminus B_t$ is a $(Q - tP)$ -negative set. In particular,

$$Q(\Omega \setminus B_0) = 0.$$

Note also that by construction,

$$B_s \supset B_t \quad \text{for } 0 \leq s < t$$

and

$$B_s = \bigcup_{t > s} B_t \quad \text{for } s \geq 0,$$

because

$$B_s = \bigcup_{\lambda \in \Lambda: \lambda > s} A_\lambda = \bigcup_{t > s} \bigcup_{\lambda \in \Lambda: \lambda > t} A_\lambda = \bigcup_{t > s} B_t.$$

The intersection

$$B_* := \bigcap_{t > 0} B_t$$

satisfies $P(B_*) = 0$, because

$$0 \leq P(B_*) \leq \inf_{t > 0} P(B_t) \leq \inf_{t > 0} Q(B_t)/t \leq \inf_{t > 0} Q(\Omega)/t = 0.$$

Next we define a function $\bar{f} : \Omega \rightarrow [0, \infty]$ via

$$\bar{f}(\omega) := \begin{cases} 0 & \text{if } \omega \in \Omega \setminus B_0, \\ \sup\{t \geq 0 : \omega \in B_t\} & \text{if } \omega \in B_0. \end{cases}$$

Note that $\{\bar{f} = \infty\} = B_*$. Note also that

$$\{\bar{f} > t\} = \bigcup_{s > t} B_s = B_t \quad \text{for all } t \geq 0,$$

so \bar{f} is measurable. Our goal is to show that $Q(A) = Q(A \cap B_*) + \int_A \bar{f} dP$ for arbitrary $A \in \mathcal{A}$. For any fixed parameter $\gamma > 1$, the set Ω may be partitioned into the disjoint sets $\Omega \setminus B_0$, B_* and

$$B_0 \setminus B_* = \bigcup_{z \in \mathbb{Z}} C_z \quad \text{with} \quad C_z := B_{\gamma^z} \setminus B_{\gamma^{z+1}}.$$

Note that $\bar{f} = 0$ on $\Omega \setminus B_0$ and $\gamma^z < \bar{f} \leq \gamma^{z+1}$ on C_z , $z \in \mathbb{Z}$. Moreover, C_z is $(Q - \gamma^z P)$ -positive and $(Q - \gamma^{z+1} P)$ -negative. Hence, for any set $A \in \mathcal{A}$,

$$Q(A \setminus B_0) = 0 = \int_{A \setminus B_0} \bar{f} dP,$$

and

$$Q(A \cap C_z) \begin{cases} \leq \gamma^{z+1} P(A \cap C_z) \leq \gamma \int_{A \cap C_z} \bar{f} dP, \\ \geq \gamma^z P(A \cap C_z) \geq \gamma^{-1} \int_{A \cap C_z} \bar{f} dP. \end{cases}$$

Consequently,

$$Q(A \setminus B_*) = \sum_{z \in \mathbb{Z}} Q(A \cap C_z) \begin{cases} \leq \gamma \sum_{z \in \mathbb{Z}} \int_{A \cap C_z} \bar{f} dP = \gamma \int_{A \setminus B_*} \bar{f} dP, \\ \geq \gamma^{-1} \sum_{z \in \mathbb{Z}} \int_{A \cap C_z} \bar{f} dP = \gamma^{-1} \int_{A \setminus B_*} \bar{f} dP, \end{cases}$$

and for $\gamma \rightarrow 1$ we obtain the desired equation(s)

$$Q(A) = Q(A \cap B_*) + \int_{A \setminus B_*} \bar{f} dP = Q(A \cap B_*) + \int_A \bar{f} dP.$$

The latter equation follows from the fact that $P(B_*) = 0$. Indeed, this representation of Q remains valid if we replace \bar{f} with the real-valued function $f := 1_{\Omega \setminus B_*} \bar{f}$. Setting $A = \Omega$ yields the equation $\int f dP = Q(\Omega) - Q(B_*) < \infty$, whence $f \in \mathcal{L}^1(P)$. \square

Theorem 2.23 shows that Q is the sum of two measures, a ‘‘singular part with respect to P ’’,

$$A \mapsto Q(A \cap B_*),$$

and an ‘‘absolutely continuous part with respect to P ’’,

$$A \mapsto \int_A f dP.$$

The singular part is nontrivial if and only if $Q(B_*) > 0$. Otherwise Q is ‘‘absolutely continuous with respect to P ’’ as defined in the next paragraph.

2.2.2 Absolute continuity and σ -finite measures

Definition 2.24 (Absolute continuity). Let P and Q be measures on (Ω, \mathcal{A}) . The measure Q is called *absolutely continuous with respect to P* if $Q(A) = 0$ for all sets $A \in \mathcal{A}$ such that $P(A) = 0$.

Definition 2.25 (σ -finiteness). A measure P on (Ω, \mathcal{A}) is called σ -finite if there exists a sequence $(A_n)_{n \geq 1}$ in \mathcal{A} such that $\Omega = \bigcup_{n \geq 1} A_n$ and $P(A_n) < \infty$ for all $n \geq 1$.

A standard example of a σ -finite measure is Lebesgue measure on \mathbb{R} . Here $\Omega = \mathbb{R}$ is the union of, say, all intervals $A_n := [-n, n]$, $n \in \mathbb{N}$.

Many results for finite measures may be extended to σ -finite measures by means of the following observation.

Lemma 2.26. For a nonzero measure P on (Ω, \mathcal{A}) , the following two statements are equivalent:

(i) P is σ -finite.

(ii) There exist a probability measure P_o on (Ω, \mathcal{A}) and an \mathcal{A} -measurable function $p : \Omega \rightarrow (0, \infty)$ such that

$$P(A) = \int_A p dP_o \quad \text{for all } A \in \mathcal{A}.$$

Proof of Lemma 2.26. Suppose first that P may be represented as in part (ii). Then $\Omega = \bigcup_{n \geq 1} A_n$ with $A_n := \{p \leq n\}$, and $P(A_n) \leq nP_o(A_n) \leq n$ for all n . Thus, P is σ -finite. Moreover, $P(\Omega) > 0$ because $\Omega = \bigcup_{m \geq 1} \{p \geq m^{-1}\}$, so $P(\{p \geq m^{-1}\}) \geq m^{-1}P_o(\{p \geq m^{-1}\}) > 0$ for sufficiently large m .

Now suppose that P is nonzero and σ -finite. In case of $P(\Omega) < \infty$, condition (ii) is satisfied with $P_o := P(\Omega)^{-1}P$ and $p \equiv P(\Omega)$. In case of $P(\Omega) = \infty$, let $\Omega = \bigcup_{n=1}^{\infty} A_n$ with sets $A_n \in \mathcal{A}$ such that $P(A_n) < \infty$, we may assume without loss of generality that these sets A_n are pairwise disjoint with $P(A_n) > 0$. But then

$$P_o(A) := \sum_{n=1}^{\infty} 2^{-n} P(A_n)^{-1} P(A \cap A_n)$$

defines a probability measure on (Ω, \mathcal{A}) , and for any $A \in \mathcal{A}$,

$$P(A) = \sum_{n=1}^{\infty} 2^n P(A_n) P_o(A \cap A_n) = \int_A p dP_o$$

with $p(\omega) := \sum_{n=1}^{\infty} 2^n P(A_n) 1_{A_n}(\omega) > 0$. □

Theorem 2.23 implies the following result about σ -finite measures.

Theorem 2.27 (Radon–Nikodym). Let P and Q be σ -finite measures on (Ω, \mathcal{A}) . Then the following two conditions are equivalent:

(i) Q is absolutely continuous with respect to P .

(ii) There exists a density f of Q with respect to P .

In case of (i–ii), the density f is P -almost everywhere unique. That means, if \tilde{f} is another density of Q with respect to P , then $P(\{\tilde{f} \neq f\}) = 0$.

Proof of Theorem 2.27. If condition (ii) is satisfied, then $P(A) = 0$ implies that $Q(A)$ equals $\int_A f dP = 0$, so condition (i) holds true as well.

Suppose that condition (i) is satisfied. If P and Q are finite, the existence of a density $f = dQ/dP$ is a consequence of Theorem 2.23. In the general case, let P_o and Q_o be probability measures on (Ω, \mathcal{A}) such that there exist strictly positive densities $p = dP/dP_o$ and $q = dQ/dQ_o$. Then $p^{-1} = dP_o/dP$ and $q^{-1} = dQ_o/dQ$, see Corollary 2.17. Hence, for arbitrary sets $A \in \mathcal{A}$,

$$\begin{aligned} P(A) &= 0 \text{ if and only if } P_o(A) = 0, \\ Q(A) &= 0 \text{ if and only if } Q_o(A) = 0. \end{aligned}$$

Thus condition (i) implies that Q_o is absolutely continuous with respect to P_o . Consequently, there exists a density $f_o = dQ_o/dP_o$. But then Lemma 2.16 implies that for any \mathcal{A} -measurable function $h : \Omega \rightarrow [0, \infty)$,

$$\int h dQ = \int h q dQ_o = \int h q f_o dP_o = \int h q f_o p^{-1} dP,$$

so $f := q f_o p^{-1} = dQ/dP$.

Suppose that \tilde{f} is another version of dQ/dP . Then, $\tilde{f}_o := q^{-1} \tilde{f} p$ is another version of dQ_o/dP_o , and $\{\tilde{f}_o \neq f_o\} = \{\tilde{f} \neq f\}$. According to Exercise 2.20, $0 = P_o\{\tilde{f}_o \neq f_o\} = P_o\{\tilde{f} \neq f\}$, whence $P\{\tilde{f} \neq f\} = \int_{\{\tilde{f} \neq f\}} p dP_o = 0$. \square

Exercise 2.28. Let P be a σ -finite measure on (Ω, \mathcal{A}) , and for some measurable function $f : \Omega \rightarrow [0, \infty)$, let $Q(A) := \int_A f dP$ for $A \in \mathcal{A}$. Show that Q is σ -finite too.

2.3 Another Riesz Representation

Combining the Radon–Nikodym theorem with Theorem 1.15 leads to a well-known representation theorem for L^p -spaces. Let $(\Omega, \mathcal{A}, \mu)$ be a σ -finite measure space. For $p \in [1, \infty)$ let $\mathcal{L}^p(\mu)$ be the set of measurable functions $f : \Omega \rightarrow \mathbb{R}$ such that

$$\|f\|_{p,\mu} := \left(\int |f|^p d\mu \right)^{1/p} < \infty.$$

It is well-known that $\mathcal{L}^p(\mu)$ is a Stone lattice and $\|\cdot\|_{p,\mu}$ defines a seminorm on $\mathcal{L}^p(\mu)$. The same is true for the set $\mathcal{L}^\infty(\mu)$ of measurable functions $f : \Omega \rightarrow \mathbb{R}$ such that

$$\|f\|_{\infty,\mu} := \inf\{r \geq 0 : \mu(\{|f| > r\}) = 0\} < \infty.$$

Theorem 2.29. For an arbitrary $p \in [1, \infty)$, let $L : \mathcal{L}^p(\mu) \rightarrow \mathbb{R}$ be linear and continuous with respect to $\|\cdot\|_{p,\mu}$. Then there exists an function $h \in \mathcal{L}^q(\mu)$, where $q = p/(p-1) \in (1, \infty]$, such that

$$L(f) = \int f h d\mu \quad \text{for all } f \in \mathcal{L}^p(\mu),$$

and

$$\sup_{f \in \mathcal{L}^p(\mu): \|f\|_{p,\mu} \leq 1} |L(f)| = \|h\|_{q,\mu}.$$

The function h is unique μ -almost everywhere.

In the proof of this theorem, the following result is useful.

Lemma 2.30 (Hölder). *Let $p \in [1, \infty)$ and $q = p/(p-1) \in (1, \infty]$. For any function $h \in \mathcal{L}^q(\mu)$,*

$$\sup_{f \in \mathcal{L}^p(\mu): \|f\|_{p,\mu} \leq 1} \left| \int f h d\mu \right| = \|h\|_{q,\mu}.$$

Proof of Lemma 2.30. If $\|h\|_{q,\mu} = 0$, the assertion is obvious, because $\mu(\{h \neq 0\}) = 0$. Thus let $0 < \|h\|_{q,\mu} < \infty$.

In case of $p = 1$, it follows from $\mu(\{|h| > \|h\|_{\infty,\mu}\}) = 0$ that

$$\left| \int f h d\mu \right| \leq \|f\|_{1,\infty} \|h\|_{\infty,\mu}.$$

On the other hand, $\mu(\{|h| \geq C\}) > 0$ for any fixed $0 < C < \|h\|_{\infty,\mu}$, and σ -finiteness of μ implies that $0 < \lambda := \mu(\{|h| \geq C\} \cap B) < \infty$ for some $B \in \mathcal{A}$. If we define $f := \lambda^{-1} \text{sign}(h) 1_{\{|h| \geq C\} \cap B}$, then $\|f\|_{1,\mu} = 1$ and

$$\int f h d\mu \geq C.$$

Hence, the supremum of $\left| \int f h d\mu \right|$ over all $f \in \mathcal{L}^1(\mu)$ with $\|f\|_{1,\mu} \leq 1$ equals $\|h\|_{\infty,\mu}$.

In case of $p > 1$, the well-known Hölder inequality states that

$$\left| \int f h d\mu \right| \leq \|f\|_{p,\mu} \|h\|_{q,\mu}$$

for arbitrary $f \in \mathcal{L}^p(\mu)$. Consequently, $\left| \int f h d\mu \right|$ is no larger than $\|h\|_{q,\mu}$ whenever $\|f\|_{p,\mu} \leq 1$. It suffices to show that we have equality for a suitable f . To this end, let $f := \lambda \text{sign}(h) |h|^{q-1} = \lambda \text{sign}(h) |h|^{1/(p-1)}$ for some $\lambda > 0$. Then $\|f\|_{p,\mu} = \lambda \|h\|_{q,\mu}^{q/p}$, so $\lambda = \|h\|_{q,\mu}^{-q/p} = \|h\|_{q,\mu}^{-1/(p-1)}$ yields a function with $\|f\|_{p,\mu} = 1$, and

$$\int f h d\mu = \|h\|_{q,\mu}^{-1/(p-1)} \int |h|^q d\mu = \|h\|_{q,\mu}.$$

□

Proof of Theorem 2.29. As in the proof of Theorem 1.15 one can show that there exist two measures M_1 and M_2 on $\sigma(\mathcal{L}^p(\mu))$ such that $\mathcal{L}^p(\mu) \subset \mathcal{L}^1(M_1 + M_2)$ and $L(f) = \int f dM_1 - \int f dM_2$ for any $f \in \mathcal{L}^p(\mu)$. Furthermore, for some constant C , $\int f dM_j \leq C \|f\|_{p,\mu}$ for all nonnegative functions $f \in \mathcal{L}^p(\mu)$ and $j = 1, 2$.

First of all, $\sigma(\mathcal{L}^p(\mu)) = \mathcal{A}$. On the one hand, $\sigma(\mathcal{L}^p(\mu)) \subset \mathcal{A}$, because all functions $f \in \mathcal{L}^p(\mu)$ are \mathcal{A} -measurable. On the other hand, if $\Omega = \bigcup_{m \geq 1} B_m$ with sets $B_m \in \mathcal{A}$ such that $\mu(B_m) < \infty$, then any set $A \in \mathcal{A}$ can be written as $A = \bigcup_{m \geq 1} A \cap B_m$, and $1_{A \cap B_m} \in \mathcal{L}^p(\mu)$ for all $m \geq 1$, whence $A \in \sigma(\mathcal{L}^p(\mu))$, that is, $\mathcal{A} \subset \sigma(\mathcal{L}^p(\mu))$.

Thus for $j = 1, 2$, M_j is a measure on \mathcal{A} such that $\int f dM_j \leq C \|f\|_{p,\mu}$ for all nonnegative $f \in \mathcal{L}^p(\mu)$. In particular, if $f = 1_A$ for some $A \in \mathcal{A}$ with $\mu(A) = 0$, then $M_j(A) = 0$. Consequently, by the Radon–Nikodym theorem, there exists a density h_j of M_j with respect to μ , and $\int f dM_j = \int f h_j d\mu \leq C \|f\|_{p,\mu}$ for all nonnegative $f \in \mathcal{L}^p(\mu)$. But this implies

that $\|h_j\|_{q,\mu} < \infty$. Indeed, let $\Omega = \bigcup_{m \geq 1} B_m$ with sets $B_1 \subset B_2 \subset B_3 \subset \dots$ such that $\mu(B_m) < \infty$. Then $h_{jm} := 1_{B_m}(h_j \wedge m) \in \mathcal{L}^q(\mu)$ satisfies $h_{jm} \uparrow h_j$ pointwise as $m \rightarrow \infty$, whence $\|h_{jm}\|_{q,\mu} \rightarrow \|h\|_{q,\mu}$ as $m \rightarrow \infty$. For any nonnegative $f \in \mathcal{L}^p(\mu)$ with $\|f\|_{p,\mu} \leq 1$,

$$C \geq \int f h_j d\mu \geq \int f h_{jm} d\mu,$$

and for a suitable such function f , the integral on the right hand side equals $\|h_{jm}\|_{q,\mu}$, see the proof of Lemma 2.30. This shows that $\|h_j\|_{q,\mu} \leq C$.

All in all, $h := h_1 - h_2$ is a function in $\mathcal{L}^q(\mu)$ such that

$$L(f) = \int f h d\mu \quad \text{for all } f \in \mathcal{L}^p(\mu),$$

and it follows from Lemma 2.30 that the supremum of $|L(f)|$ over all $f \in \mathcal{L}^p(\mu)$ with $\|f\|_{p,\mu} \leq 1$ is equal to $\|h\|_{q,\mu}$.

Essential uniqueness of h can be verified as follows: If $h, \tilde{h} \in \mathcal{L}^q(\mu)$ such that $L(f) = \int f h d\mu = \int f \tilde{h} d\mu$ for all $f \in \mathcal{L}^p(\mu)$, then $\int f(h - \tilde{h}) d\mu = 0$ for all $f \in \mathcal{L}^p(\mu)$. Then it follows from Lemma 2.30 that $\|h - \tilde{h}\|_{q,\mu} = 0$, that is, $\mu(\{h \neq \tilde{h}\}) = 0$. \square

Chapter 3

Conditional Expectations

Throughout this chapter let (Ω, \mathcal{A}, P) be a probability space, and let X be a random variable in $\mathcal{L}^1(P)$. That means, $X : \Omega \rightarrow \mathbb{R}$ is \mathcal{A} -measurable, and $\int |X| dP < \infty$.

3.1 Conditional expectations with respect to a sub- σ -field

Let \mathcal{A}_o be a sub- σ -field of \mathcal{A} . One could think about a rather complex random experiment described by (Ω, \mathcal{A}, P) , and the subfield \mathcal{A}_o represents some partial aspects of it.

Theorem 3.1. *There exists a random variable $X_o \in \mathcal{L}^1(P|_{\mathcal{A}_o})$ such that*

$$(3.1) \quad \int_{A_o} X_o dP = \int_{A_o} X dP \quad \text{for all } A_o \in \mathcal{A}_o.$$

This random variable X_o is almost everywhere unique: If \tilde{X}_o is another random variable satisfying (3.1), then $P(\tilde{X}_o \neq X_o) = 0$.

Definition 3.2 (Conditional expectation, I). A random variable X_o as in Theorem 3.1 is called (a version of the) *conditional expectation of X , given the sub- σ -field \mathcal{A}_o* . A function X_o satisfying (3.1) is denoted by $\mathbb{E}(X | \mathcal{A}_o)$.

Proof of Theorem 3.1. Uniqueness of X_o almost everywhere follows from Lemma 2.20 (with \mathcal{A}_o in place of \mathcal{A}), so it suffices to prove existence of X_o .

Let $X = X^+ - X^-$ with $X^\pm = \max(\pm X, 0)$. Then

$$Q^\pm(A) := \int_A X^\pm dP$$

defines finite measures on (Ω, \mathcal{A}) with density X^\pm with respect to P . This implies that $Q^\pm|_{\mathcal{A}_o}$ is absolutely continuous with respect to $P|_{\mathcal{A}_o}$. Hence, by the Radon–Nikodym theorem, there exist densities $X_o^{(\pm)} \in \mathcal{L}^1(P|_{\mathcal{A}_o})$ of $Q^\pm|_{\mathcal{A}_o}$ with respect to $P|_{\mathcal{A}_o}$. This implies that $X_o := X_o^{(+)} - X_o^{(-)}$ has the desired property (3.1). \square

Example 3.3 (Countable partitions). Let $\Omega = \bigcup_{n \geq 1} B_n$ with a sequence $(B_n)_{n \geq 1}$ of disjoint measurable sets, and let \mathcal{A}_o be the smallest σ -field containing all these sets B_n , $n \geq 1$. Then

$$\mathbb{E}(X | \mathcal{A}_o) = \sum_{n \geq 1} \mathbb{E}(X | B_n) \cdot 1_{B_n} \quad \text{almost surely}$$

with the numbers

$$\mathbb{E}(X | B) := \begin{cases} P(B)^{-1} \int_B X dP & \text{if } P(B) > 0, \\ \mathbb{E}(X) & \text{else.} \end{cases}$$

(The definition $\mathbb{E}(X | B) := \mathbb{E}(X)$ in case of $P(B) = 0$ is somewhat arbitrary.) Indeed, the function $X_o := \sum_{n \geq 1} \mathbb{E}(X | B_n) \cdot 1_{B_n}$ is constant on each set B_n , $n \geq 1$, so it is \mathcal{A}_o -measurable. Moreover, any set $A_o \in \mathcal{A}_o$ may be written as $A_o = \bigcup_{n \in M} B_n$ with $M \subset \mathbb{N}$, so

$$\int_{A_o} X dP = \sum_{n \in M} \int_{B_n} X dP = \sum_{n \in M} P(B_n) \mathbb{E}(X | B_n) = \sum_{n \in M} \int_{B_n} X_o dP = \int_{A_o} X_o dP.$$

Remark 3.4 (Properties of conditional expectations). In what follows, let $X, Y \in \mathcal{L}^1(P)$.

(a) For real numbers a, b ,

$$\mathbb{E}(aX + bY | \mathcal{A}_o) = a \mathbb{E}(X | \mathcal{A}_o) + b \mathbb{E}(Y | \mathcal{A}_o) \quad \text{almost surely.}$$

This is a simple consequence of linearity of integrals.

(b) If $X \leq Y$ almost surely, then

$$\mathbb{E}(X | \mathcal{A}_o) \leq \mathbb{E}(Y | \mathcal{A}_o) \quad \text{almost surely.}$$

This is essentially a consequence of Lemma 2.20: It follows from $X \leq Y$ almost surely that

$$\int_{A_o} \mathbb{E}(X | \mathcal{A}_o) dP = \int_{A_o} X dP \leq \int_{A_o} Y dP = \int_{A_o} \mathbb{E}(Y | \mathcal{A}_o) dP \quad \text{for all } A_o \in \mathcal{A}_o.$$

Hence, $P(\mathbb{E}(X | \mathcal{A}_o) > \mathbb{E}(Y | \mathcal{A}_o)) = 0$.

(c) The mapping $X \mapsto \mathbb{E}(X | \mathcal{A}_o)$ is a weak contraction in the sense that

$$|\mathbb{E}(X | \mathcal{A}_o)| \leq \mathbb{E}(|X| | \mathcal{A}_o) \quad \text{almost surely.}$$

This follows from properties (a–b): Writing $X = X^+ - X^-$ with $X^\pm = \max(\pm X, 0)$, one may write $|X| = X^+ + X^-$, and we know that almost everywhere,

$$\mathbb{E}(X^\pm | \mathcal{A}_o) \geq 0 \quad \text{(property (b))},$$

$$\mathbb{E}(X | \mathcal{A}_o) = \mathbb{E}(X^+ - X^- | \mathcal{A}_o) = \mathbb{E}(X^+ | \mathcal{A}_o) - \mathbb{E}(X^- | \mathcal{A}_o) \quad \text{(property (a))},$$

$$\mathbb{E}(|X| | \mathcal{A}_o) = \mathbb{E}(X^+ + X^- | \mathcal{A}_o) = \mathbb{E}(X^+ | \mathcal{A}_o) + \mathbb{E}(X^- | \mathcal{A}_o) \quad \text{(property (a))},$$

whence

$$|\mathbb{E}(X | \mathcal{A}_o)| \leq \mathbb{E}(|X| | \mathcal{A}_o).$$

(d) It follows from properties (a) and (c) that

$$\int |\mathbb{E}(X | \mathcal{A}_o) - \mathbb{E}(Y | \mathcal{A}_o)| dP \leq \int |X - Y| dP.$$

(e) For any bounded, \mathcal{A}_o -measurable random variable $Z_o : \Omega \rightarrow \mathbb{R}$,

$$(3.2) \quad \int X Z_o dP = \int \mathbb{E}(X | \mathcal{A}_o) Z_o dP.$$

Indeed, by definition of $\mathbb{E}(X | \mathcal{A}_o)$, equation (3.2) is true for indicator functions $Z_o = 1_{A_o}$, $A_o \in \mathcal{A}_o$. By linearity of integrals, (3.2) is even true for \mathcal{A}_o -measurable functions $Z_o : \Omega \rightarrow \mathbb{R}$ taking only finitely many values. But for any bounded, \mathcal{A}_o -measurable function $Z_o : \Omega \rightarrow \mathbb{R}$, there exists a sequence $(Z_n)_{n \geq 1}$ of \mathcal{A}_o -measurable functions taking only finitely many values such that

$$\delta_n := \sup_{\omega \in \Omega} |Z_n(\omega) - Z_o(\omega)| \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Now (3.2) follows from the fact that for arbitrary $n \geq 1$,

$$\begin{aligned} \left| \int X Z_o dP - \int \mathbb{E}(X | \mathcal{A}_o) Z_o dP \right| &= \left| \int (X - \mathbb{E}(X | \mathcal{A}_o))(Z_o - Z_n) dP \right| \\ &\leq \delta_n \int |X - \mathbb{E}(X | \mathcal{A}_o)| dP. \end{aligned}$$

Exercise 3.5. Let P be the exponential distribution on $\Omega = [0, \infty)$ with rate parameter $\lambda > 0$, and let \mathcal{A}_o be the smallest σ -field containing all intervals $[k, k+1)$, $k \in \mathbb{N}_0$. Determine $\mathbb{E}(X | \mathcal{A}_o)$ in case of $X(\omega) := \omega$.

Exercise 3.6 (“Tower property” of conditional expectations). Let (Ω, \mathcal{A}, P) be a probability space and $X \in \mathcal{L}^1(P)$. Let \mathcal{A}_1 and \mathcal{A}_2 be σ -fields over Ω such that $\mathcal{A}_1 \subset \mathcal{A}_2 \subset \mathcal{A}$.

(a) Show that $\mathbb{E}(\mathbb{E}(X | \mathcal{A}_1) | \mathcal{A}_2) = \mathbb{E}(X | \mathcal{A}_1)$ almost surely.

(b) Show the so-called tower property: $\mathbb{E}(\mathbb{E}(X | \mathcal{A}_2) | \mathcal{A}_1) = \mathbb{E}(X | \mathcal{A}_1)$ almost surely.

We end this section with an extension of Jensen’s inequality for expectations to conditional expectations.

Lemma 3.7 (Jensen’s inequality). Let $X \in \mathcal{L}^1(P)$, and let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be convex such that $\int \psi(X) dP < \infty$. Then

$$\psi(\mathbb{E}(X | \mathcal{A}_o)) \leq \mathbb{E}(\psi(X) | \mathcal{A}_o) \quad \text{almost surely,}$$

and

$$\int \psi(\mathbb{E}(X | \mathcal{A}_o)) dP \leq \int \psi(X) dP.$$

Proof of Lemma 3.7. By the definition of conditional expectations and Lemma 2.20, it suffices to show that for any fixed $A_o \in \mathcal{A}_o$,

$$\int_{A_o} \psi(\mathbb{E}(X | \mathcal{A}_o)) dP \leq \int_{A_o} \psi(X) dP.$$

Note that by convexity of ψ , for any fixed $x_o \in \mathbb{R}$ there exists a slope $b(x_o) \in \mathbb{R}$ such that $\psi(x) \geq \psi(x_o) + b(x_o)(x - x_o)$ for all $x \in \mathbb{R}$; see, for instance, Chapter 3 of Dümbgen (2021). Moreover, for arbitrary $x \in \mathbb{R}$,

$$\psi(x) = \sup_{x_o \in \mathbb{Q}} (\psi(x_o) + b(x_o)(x - x_o)).$$

Consequently,

$$\begin{aligned} \int_{A_o} \psi(\mathbb{E}(X | \mathcal{A}_o)) dP &= \int_{A_o} \sup_{x_o \in \mathbb{Q}} (\psi(x_o) + b(x_o)(\mathbb{E}(X | \mathcal{A}_o) - x)) dP \\ &= \int_{A_o} \sup_{x_o \in \mathbb{Q}} \mathbb{E}(\psi(x_o) + b(x_o)(X - x_o) | \mathcal{A}_o) dP \\ &\leq \int_{A_o} \psi(X) dP, \end{aligned}$$

where the second last step follows from Remark 3.4 (a) and countability of \mathbb{Q} , while the last step follows from Remark 3.4 (b) and countability of \mathbb{Q} . \square

3.2 Conditional expectations as orthogonal projections

In this section we restrict our attention to the space $\mathcal{L}^2(P) \subset \mathcal{L}^1(P)$ of square-integrable random variables. If we identify random variables which are equal almost surely, we obtain the Hilbert space $L^2(P)$ with inner product

$$\langle X, Y \rangle := \int XY dP$$

and norm

$$\|X\| := \langle X, X \rangle^{1/2} = \left(\int X^2 dP \right)^{1/2}.$$

One may view $\mathbb{E}(\cdot | \mathcal{A}_o)$ as a continuous linear mapping from $L^2(P)$ to its closed linear subspace $L^2(P|_{\mathcal{A}_o})$. Indeed, applying Lemma 3.7 with $\psi(x) := x^2$ leads to the inequality

$$\int \mathbb{E}(X | \mathcal{A}_o)^2 dP \leq \int X^2 dP.$$

Theorem 3.8. *The mapping $X \mapsto \mathbb{E}(X | \mathcal{A}_o)$ is the orthogonal linear projection of $L^2(P)$ onto $L^2(P|_{\mathcal{A}_o})$. In particular, for arbitrary random variables $X \in L^2(P)$ and $Y_o \in L^2(P|_{\mathcal{A}_o})$,*

$$\|X - Y_o\|^2 = \|X - \mathbb{E}(X | \mathcal{A}_o)\|^2 + \|\mathbb{E}(X | \mathcal{A}_o) - Y_o\|^2 \geq \|X - \mathbb{E}(X | \mathcal{A}_o)\|^2$$

with equality if and only if $Y_o = \mathbb{E}(X | \mathcal{A}_o)$. Moreover,

$$\|X\|^2 = \|X - \mathbb{E}(X | \mathcal{A}_o)\|^2 + \|\mathbb{E}(X | \mathcal{A}_o)\|^2 \geq \|\mathbb{E}(X | \mathcal{A}_o)\|^2$$

with equality if and only if $X = \mathbb{E}(X | \mathcal{A}_o)$.

Proof of Theorem 3.8. It is well-known from general results about Hilbert spaces that $X_o := \mathbb{E}(X | \mathcal{A}_o)$ is the orthogonal projection of X onto $L^2(P|_{\mathcal{A}_o})$ if and only if

$$X - X_o \perp L^2(P|_{\mathcal{A}_o});$$

see, for instance, Chapter 2 of Dümbgen (2021). That means,

$$(3.3) \quad \langle X - X_o, Z_o \rangle = 0 \quad \text{for all } Z_o \in L^2(P|_{\mathcal{A}_o}).$$

Note first, that the left hand side of (3.3) equals

$$\int X Z_o dP - \int \mathbb{E}(X | \mathcal{A}_o) Z_o dP.$$

This is equal to 0 whenever Z_o is bounded, see Remark 3.4 (e). But for arbitrary $Z_o \in L^2(P|_{\mathcal{A}_o})$,

$$Z_n := \text{sign}(Z_o) \min(|Z_o|, n)$$

defines a sequence $(Z_n)_{n \geq 1}$ of bounded random variables $Z_n \in L^2(P|_{\mathcal{A}_o})$ such that

$$\|Z_o - Z_n\| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

which may be verified by dominated convergence. Now (3.3) follows from the fact that for arbitrary $n \geq 1$,

$$|\langle X - X_o, Z_o \rangle| = |\langle X - X_o, Z_o - Z_n \rangle| \leq \|X - X_o\| \|Z_o - Z_n\|$$

by the Cauchy–Schwarz inequality.

Orthogonality of $X - X_o$ onto $L^2(P|_{\mathcal{A}_o})$ implies that for any $Y_o \in L^2(P|_{\mathcal{A}_o})$,

$$\begin{aligned} \|X - Y_o\|^2 &= \|X - X_o + X_o - Y_o\|^2 \\ &= \|X - X_o\|^2 + \|X_o - Y_o\|^2 + 2\langle X - X_o, X_o - Y_o \rangle \\ &= \|X - X_o\|^2 + \|X_o - Y_o\|^2, \end{aligned}$$

and in the special case $Y_o = 0$ we obtain

$$\|X\|^2 = \|X - X_o\|^2 + \|X_o\|^2.$$

□

Exercise 3.9. As in Exercise 3.5 let P be the exponential distribution on $\Omega = [0, \infty)$ with rate parameter $\lambda > 0$, let \mathcal{A}_o be the smallest σ -field containing all intervals $[k, k + 1)$, $k \in \mathbb{N}_0$, and let $X(\omega) := \omega$. Determine the value of

$$\|X - \mathbb{E}(X | \mathcal{A}_o)\|.$$

Hint for checking your solution: $\|X - \mathbb{E}(X | \mathcal{A}_o)\|^2$ equals $\lambda^{-2} - e^\lambda(e^\lambda - 1)^{-2}$.

3.3 Conditional expectations given another random variable

Previously we referred to \mathcal{A}_o as describing partial aspects of the random experiment (Ω, \mathcal{A}, P) . Specifically, suppose that instead of the random outcome $\omega \in \Omega$ we only know $T(\omega)$ for some given measurable function $T : (\Omega, \mathcal{A}) \rightarrow (\mathcal{T}, \mathcal{B})$. This corresponds to the sub- σ -field

$$\sigma(T) := \{T^{-1}(B) : B \in \mathcal{B}\}$$

of \mathcal{A} . It is the smallest σ -field \mathcal{A}_o over Ω such that T is \mathcal{A}_o - \mathcal{B} -measurable. There is a simple result about $\sigma(T)$ -measurable functions $X : \Omega \rightarrow \mathbb{R}$.

Lemma 3.10 (Lifting). *A mapping $X : \Omega \rightarrow \mathbb{R}$ is $\sigma(T)$ -measurable if and only if*

$$X = V \circ T$$

for some \mathcal{B} -measurable function $V : \mathcal{T} \rightarrow \mathbb{R}$.

Here $V \circ T$ stands for the mapping $\Omega \ni \omega \mapsto V(T(\omega)) \in \mathbb{R}$. In connection with integrals we often write $V(T)$ instead of $V \circ T$.

Proof of Lemma 3.10. On the one hand, if $X = V \circ T$ with a \mathcal{B} -measurable function $V : \mathcal{T} \rightarrow \mathbb{R}$, then $V^{-1}(C) \in \mathcal{B}$ and $X^{-1}(C) = T^{-1}(V^{-1}(C)) \in \sigma(T)$ for any Borel set $C \subset \mathbb{R}$.

On the other hand, let X be $\sigma(T)$ -measurable. For $n \in \mathbb{N}$ let $X_n := 2^{-n} \lfloor 2^n X \rfloor$. That means, for $z \in \mathbb{Z}$,

$$A_{n,z} := \{X_n = 2^{-n}z\} = \{2^{-n}z \leq X < 2^{-n}(z+1)\}.$$

By assumption on X , the latter event may be written as $\{T \in B_{n,z}\}$ for some set $B_{n,z} \in \mathcal{B}$. Since the sets $A_{n,z}$, $z \in \mathbb{Z}$, are pairwise disjoint, we may assume without loss of generality that the sets $B_{n,z}$, $z \in \mathbb{Z}$, are pairwise disjoint too. Just use some enumeration $z(1), z(2), z(3), \dots$ of \mathbb{Z} , and then replace $B_{n,z(k)}$ with its subset $B_{n,z(k)} \setminus \bigcup_{\ell < k} B_{n,z(\ell)}$. Then,

$$X_n = V_n \circ T$$

with the \mathcal{B} -measurable function $V_n : \mathcal{T} \rightarrow \mathbb{R}$,

$$V_n(t) := \begin{cases} 2^{-n}z & \text{if } t \in B_{n,z}, z \in \mathbb{Z}, \\ 0 & \text{if } t \notin \bigcup_{z \in \mathbb{Z}} B_{n,z}. \end{cases}$$

Since $X_n \leq X < X_n + 2^{-n}$,

$$X = \lim_{n \rightarrow \infty} V_n \circ T = \bar{V} \circ T$$

with the \mathcal{B} -measurable function $\bar{V} : \mathcal{T} \rightarrow \bar{\mathbb{R}}$,

$$\bar{V}(t) := \limsup_{n \rightarrow \infty} V_n(t).$$

But $X = \bar{V} \circ T$ being real-valued implies that $T^{-1}(\bar{V}^{-1}(\{-\infty, \infty\})) = \emptyset$. Hence, $X = V \circ T$ with the \mathcal{B} -measurable function $V : \mathcal{T} \rightarrow \mathbb{R}$ given by, say,

$$V(t) := \begin{cases} \bar{V}(t) & \text{if } \bar{V}(t) \in \mathbb{R}, \\ 0 & \text{else.} \end{cases}$$

□

In addition to Lemma 3.10 we need an elementary result about the distribution P^T of T , i.e. the probability measure

$$\mathcal{B} \ni B \mapsto P^T(B) := P(T \in B).$$

Exercise 3.11 (Change of variables). Show that for any real-valued random variable V on the probability space $(\mathcal{T}, \mathcal{B}, P^T)$ and any set $B \in \mathcal{B}$,

$$\int_{T^{-1}(B)} V(T) dP = \int_B V dP^T,$$

provided that one of these two integrals is well-defined in $\overline{\mathbb{R}}$. Hint: Use approximations of V as in Lemma 2.16.

Now we can prove a generalization of Theorem 3.1:

Theorem 3.12. For any random variable $X \in \mathcal{L}^1(P)$ there exists a random variable $V \in \mathcal{L}^1(P^T)$ such that

$$(3.4) \quad \int_{T^{-1}(B)} X dP = \int_B V dP^T \quad \text{for arbitrary } B \in \mathcal{B}.$$

This random variable V is almost everywhere unique in the following sense: If \tilde{V} is another random variable satisfying (3.4), then $P^T(\tilde{V} \neq V) = 0$.

Definition 3.13 (Conditional expectation, II). A random variable V as in Theorem 3.12 is called (a version of the) conditional expectation of X , given the random variable T . For any such random variable V we write $\mathbb{E}(X | T)$ instead of V and $\mathbb{E}(X | T = t)$ instead of $\mathbb{E}(V | T)(t)$ or $V(t)$.

Note that $\mathbb{E}(X | T)$ is a function on \mathcal{T} . But sometimes we and other people abuse notation slightly to write $\mathbb{E}(X | T)$ instead of $\mathbb{E}(X | \sigma(T))$, which is a function on Ω .

Example 3.14. Let $\Omega = [0, \infty)$, and let P be the exponential distribution with rate parameter $\lambda > 0$, i.e. with density $f_\lambda(\omega) := \lambda e^{-\lambda\omega}$ with respect to Lebesgue measure on Ω . Let $X(\omega) := \omega$, and let \mathcal{A}_o be the smallest σ -field over Ω containing all intervals $[k, k + 1)$, $k \in \mathbb{N}_0$. One can easily verify that $\mathcal{A}_o = \sigma(T)$ with

$$T : [0, \infty) \rightarrow \mathbb{N}_0, \quad T(\omega) := \lfloor \omega \rfloor.$$

And the result of Exercise 3.5 may be reformulated as

$$\begin{aligned} \mathbb{E}(X | \mathcal{A}_o) &= \lfloor X \rfloor + a(\lambda), \\ \mathbb{E}(X | T = t) &= t + a(\lambda) \quad \text{for } t \in \mathbb{N}_0, \end{aligned}$$

with $a(\lambda) = \lambda^{-1}(e^\lambda - 1 - \lambda)/(e^\lambda - 1)$.

Remark 3.15. Theorem 3.1 may be viewed as a special case of Theorem 3.12 if we define $(\mathcal{T}, \mathcal{B}) := (\Omega, \mathcal{A}_o)$ and $T(\omega) := \omega$.

Proof of Theorem 3.12. We may apply Theorem 3.1 to the sub- σ -field $\mathcal{A}_o = \sigma(T)$. This yields a random variable $X_o \in \mathcal{L}^1(P|_{\sigma(T)})$ such that

$$\int_{T^{-1}(B)} X dP = \int_{T^{-1}(B)} X_o dP \quad \text{for all } B \in \mathcal{B}.$$

Lemma 3.10 implies that $X_o = V \circ T$ for some \mathcal{B} -measurable function $V : \mathcal{T} \rightarrow \mathbb{R}$. Now it follows from Exercise 3.11 that

$$\int_{T^{-1}(B)} V(T) dP = \int_B V dP^T \quad \text{for all } B \in \mathcal{B}.$$

This proves existence of a random variable V with the desired property (3.4). Essential uniqueness of V is a consequence of Lemma 2.20. \square

Remark 3.16 (Properties of conditional expectations). The properties listed in Remark 3.4 carry over to the present setting with only few modifications. In what follows, let $X, Y \in \mathcal{L}^1(P)$.

(a) For real numbers a, b ,

$$\mathbb{E}(aX + bY | T) = a \mathbb{E}(X | T) + b \mathbb{E}(Y | T) \quad P^T\text{-almost surely.}$$

(b) If $X \leq Y$ almost surely, then

$$\mathbb{E}(X | T) \leq \mathbb{E}(Y | T) \quad P^T\text{-almost surely.}$$

(c) The mapping $X \mapsto \mathbb{E}(X | T)$ is a weak contraction in the sense that

$$|\mathbb{E}(X | T)| \leq \mathbb{E}(|X| | T) \quad P^T\text{-almost surely.}$$

(d) It follows from properties (a) and (c) that

$$\int |\mathbb{E}(X | T) - \mathbb{E}(Y | T)| dP^T \leq \int |X - Y| dP.$$

(e) For any bounded, \mathcal{B} -measurable function $h : \mathcal{T} \rightarrow \mathbb{R}$,

$$(3.5) \quad \int Xh(T) dP = \int \mathbb{E}(X | T)h dP^T.$$

This follows again from an approximation argument. The definition of $\mathbb{E}(X | T)$ implies that equation (3.5) is true for indicator functions $h = 1_B$, $B \in \mathcal{B}$, because $1_{T^{-1}(B)} = 1_B(T)$. The remaining arguments are analogous to the arguments for Remark 3.4 (e).

An extension of the lifting lemma. Although it is not needed in this lecture, it is worthwhile to mention that Lemma 3.10 can be generalised as follows:

Lemma 3.17 (Lifting). *Let (\mathcal{X}, d) be a complete, separable metric space, equipped with its Borel- σ -field. A mapping $X : \Omega \rightarrow \mathcal{X}$ is $\sigma(T)$ -measurable if and only if*

$$X = V \circ T$$

for some \mathcal{B} -measurable function $V : \mathcal{T} \rightarrow \mathcal{X}$.

The proof is split into the following two exercises.

Exercise 3.18. Let (Ω, \mathcal{A}) be a measurable and (\mathcal{X}, d) a metric space. Further let $(V_n)_{n \geq 1}$ be a sequence of measurable functions $V_n : \Omega \rightarrow \mathcal{X}$, where \mathcal{X} is equipped with its Borel- σ -field with respect to d .

(a) Suppose that $(V_n)_{n \geq 1}$ converges pointwise to a function $V : \Omega \rightarrow \mathcal{X}$. Show that V is measurable.

(b) Suppose that (\mathcal{X}, d) is separable and complete. Without any further assumptions on $(V_n)_{n \geq 1}$, show that the set C of all $\omega \in \Omega$ such that $(V_n(\omega))_{n \geq 1}$ converges is a measurable subset of Ω .

Exercise 3.19. Imitate the proof of Lemma 3.10 to construct a sequence $(V_n)_{n \geq 1}$ of measurable functions $V_n : \mathcal{T} \rightarrow \mathcal{X}$ converging pointwise to a function V such that $X = V \circ T$. Conclude from Exercise 3.18 that V is measurable.

Chapter 4

Stochastic Kernels

In this section we collect some useful results about finite measures on product spaces. Throughout let $(\mathcal{X}, \mathcal{A})$ and $(\mathcal{Y}, \mathcal{B})$ be measurable spaces, and we consider the Cartesian product $\Omega := \mathcal{X} \times \mathcal{Y}$ equipped with the product σ -field

$$\mathcal{C} := \mathcal{A} \otimes \mathcal{B},$$

i.e. the smallest σ -field containing all sets $A \times B$ with $A \in \mathcal{A}$ and $B \in \mathcal{B}$. An important property of \mathcal{C} is that for any $C \in \mathcal{C}$,

$$(4.1) \quad \{x \in \mathcal{X} : (x, y_o) \in C\} \in \mathcal{A} \quad \text{for all } y_o \in \mathcal{Y},$$

$$(4.2) \quad \{y \in \mathcal{Y} : (x_o, y) \in C\} \in \mathcal{B} \quad \text{for all } x_o \in \mathcal{X}.$$

The reason is that the family of all sets $C \subset \Omega$ satisfying (4.1) and (4.2) is easily verified to be a σ -field over Ω containing the generator

$$\mathcal{A} \boxtimes \mathcal{B} := \{A \times B : A \in \mathcal{A}, B \in \mathcal{B}\}$$

of \mathcal{C} .

4.1 Stochastic Kernels and Fubini's Theorem

Definition 4.1 (Stochastic kernel). A stochastic kernel from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ is a mapping

$$K : \mathcal{X} \times \mathcal{B} \rightarrow [0, 1]$$

such that

- for any fixed $x \in \mathcal{X}$, the mapping $K(x, \cdot)$ defines a probability measure on $(\mathcal{Y}, \mathcal{B})$,
- for any fixed $B \in \mathcal{B}$, the mapping $K(\cdot, B)$ is \mathcal{A} -measurable on \mathcal{X} .

Remark 4.2 (Randomized mappings). One may interpret a stochastic kernel K from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ as a “randomized measurable mapping” from \mathcal{X} to \mathcal{Y} . Instead of mapping a point $x \in \mathcal{X}$ to a unique point $k(x) \in \mathcal{Y}$, we choose a random point in \mathcal{Y} with distribution $K(x, \cdot)$. Indeed, any \mathcal{A} - \mathcal{B} -measurable mapping $k : \mathcal{X} \rightarrow \mathcal{Y}$ corresponds to the (degenerate) stochastic kernel K given by

$$K(x, \cdot) := \delta_{k(x)}, \quad K(x, B) = 1_{k^{-1}(B)}(x).$$

Theorem 4.3 (Product of a finite measure and a stochastic kernel). *Let P be a finite measure on $(\mathcal{X}, \mathcal{A})$, and let K be a stochastic kernel from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$. Then for arbitrary $C \in \mathcal{C}$, the number*

$$P \otimes K(C) := \int_{\mathcal{X}} K(x, C_x) P(dx) \quad \text{with } C_x := \{y \in \mathcal{Y} : (x, y) \in C\}$$

is well-defined, and $P \otimes K$ is a finite measure on (Ω, \mathcal{C}) with $P \otimes K(\Omega) = P(\mathcal{X})$.

Remark 4.4 (Product measures). Let Q be a probability measure on $(\mathcal{Y}, \mathcal{B})$. If we set $K(x, B) := Q(B)$ for arbitrary $x \in \mathcal{X}$ and $B \in \mathcal{B}$, then $P \otimes K$ is the usual product measure $P \otimes Q$.

In the subsequent proofs, we use repeatedly the notion of a Dynkin system, see Section A.1 for the basic concepts.

Proof of Theorem 4.3. One can easily verify that the family \mathcal{D} of sets $C \in \mathcal{C}$ such that $x \mapsto K(x, C_x)$ is \mathcal{A} -measurable is a Dynkin system containing the generator $\mathcal{A} \boxtimes \mathcal{B}$ of \mathcal{C} . Since $\mathcal{A} \boxtimes \mathcal{B}$ is closed under intersections, the smallest Dynkin system containing $\mathcal{A} \boxtimes \mathcal{B}$ coincides with the product σ -field $\mathcal{A} \otimes \mathcal{B} = \mathcal{C}$. Hence $P \otimes K(C)$ is well-defined for any $C \in \mathcal{C}$.

It follows from the properties of K that $P \otimes K$ defines a content, i.e. a finitely additive function on $\mathcal{A} \otimes \mathcal{B}$ with $P \otimes K(\emptyset) = 0$ and $P \otimes K(\Omega) = P(\mathcal{X})$. It is a measure, because for arbitrary sets $C^{(1)} \subset C^{(2)} \subset C^{(3)} \subset \dots$ in \mathcal{C} and $C := \bigcup_{n \geq 1} C^{(n)}$,

$$\begin{aligned} \lim_{n \rightarrow \infty} P \otimes K(C^{(n)}) &= \lim_{n \rightarrow \infty} \int K(x, C_x^{(n)}) P(dx) \\ &= \int \lim_{n \rightarrow \infty} K(x, C_x^{(n)}) P(dx) \\ &= \int K(x, C_x) P(dx) = P \otimes K(C) \end{aligned}$$

by monotone convergence. □

If a distribution on a product space is a product of a marginal distribution and a stochastic kernel, there is a generalization of Fubini's theorem for integrals with respect to product measures.

Theorem 4.5 (Fubini's theorem). *Let $\mathbb{P} = P \otimes K$ with a finite measure P on $(\mathcal{X}, \mathcal{A})$ and a stochastic kernel K from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$.*

(i) *For any \mathcal{C} -measurable function $f : \Omega \rightarrow [0, \infty]$,*

$$h(x) := \int_{\mathcal{Y}} f(x, y) K(x, dy)$$

defines an \mathcal{A} -measurable mapping $h : \mathcal{X} \rightarrow [0, \infty]$, and

$$\int_{\Omega} f d\mathbb{P} = \int_{\mathcal{X}} h dP.$$

(ii) *If $f \in \mathcal{L}^1(\mathbb{P})$, then there exists a set $A(f) \in \mathcal{A}$ with $P(A(f)) = 0$ such that for all $x \in \mathcal{X} \setminus A(f)$,*

$$h(x) := \int f(x, y) K(x, dy)$$

is well-defined in \mathbb{R} , and

$$\int_{\Omega} f d\mathbb{P} = \int_{\mathcal{X}} h dP.$$

Proof of Theorem 4.5. Part (ii) is an immediate consequence of part (i). Hence we prove only part (i). For a measurable function $f : \Omega \rightarrow [0, \infty]$ and $n \in \mathbb{N}$ let

$$f_n := 2^{-n} \sum_{k=1}^{n2^n} 1_{[f \geq k2^{-n}]}.$$

This defines a sequence of measurable functions $f_n \geq 0$ such that $f_n \uparrow f$ as $n \rightarrow \infty$ pointwise. Indeed, for any integer $\ell \geq 0$,

$$f_n = \min(\ell, n2^n)2^{-n} \quad \text{on} \quad \{\ell 2^{-n} \leq f < (\ell + 1)2^{-n}\}.$$

For any fixed n ,

$$\begin{aligned} \int_{\Omega} f_n d\mathbb{P} &= 2^{-n} \sum_{k=1}^{n2^n} \mathbb{P}(f \geq k2^{-n}) \\ &= 2^{-n} \sum_{k=1}^{n2^n} \int_{\mathcal{X}} K(x, \{f \geq k2^{-n}\}_x) P(dx) \\ &= \int_{\mathcal{X}} h_n dP \end{aligned}$$

with

$$h_n(x) := 2^{-n} \sum_{k=1}^{n2^n} K(x, \{f \geq k2^{-n}\}_x) = \int_{\mathcal{Y}} f_n(x, y) K(x, dy).$$

This is obviously a measurable function on $(\mathcal{X}, \mathcal{A})$. By monotone convergence, for any $x \in \mathcal{X}$,

$$h_n(x) \uparrow h(x) = \int_{\mathcal{Y}} f(x, y) K(x, dy) \quad \text{as } n \rightarrow \infty,$$

whence h is measurable as well. Another application of monotone convergence yields that

$$\int_{\mathcal{X}} h_n dP \uparrow \int_{\mathcal{X}} h dP \quad \text{as } n \rightarrow \infty.$$

On the other hand, by monotone convergence,

$$\int_{\Omega} f_n d\mathbb{P} \uparrow \int_{\Omega} f d\mathbb{P} \quad \text{as } n \rightarrow \infty,$$

so the asserted equation is true. \square

Remark 4.6 (Fubini's theorem for products of σ -finite measures). The present results also imply a more traditional version of Fubini's theorem. Let P and Q be σ -finite measures on $(\mathcal{X}, \mathcal{A})$ and on $(\mathcal{Y}, \mathcal{B})$, respectively. For $C \in \mathcal{C}$, the number

$$(4.3) \quad \mathbb{P}(C) := \int_{\mathcal{X}} Q(C_x) P(dx) \quad \text{with } C_x := \{y \in \mathcal{Y} : (x, y) \in C\}$$

is well-defined in $[0, \infty]$ and defines a σ -finite measure on (Ω, \mathcal{C}) . Moreover, for arbitrary \mathcal{C} -measurable functions $f : \Omega \rightarrow [0, \infty]$,

$$h(x) := \int_{\mathcal{Y}} f(x, y) Q(dy)$$

defines a \mathcal{A} -measurable function $h : \mathcal{X} \rightarrow [0, \infty]$, and

$$\int_{\Omega} f d\mathbb{P} = \int_{\mathcal{X}} h dP.$$

To verify this, we apply Lemma 2.26 to find probability measures P_o on $(\mathcal{X}, \mathcal{A})$ and Q_o on $(\mathcal{Y}, \mathcal{B})$ such that densities $p = dP/dP_o$ and $q = dQ/dQ_o$ exist. Then the product measure $\mathbb{P}_o := P_o \otimes Q_o$ is well-defined, and

$$\mathbb{P}(C) := \int_C p(x)q(y) \mathbb{P}_o(d(x, y))$$

defines a σ -finite measure on (Ω, \mathcal{C}) . This definition coincides with (4.3), because Theorem 4.5 implies that

$$\begin{aligned} \mathbb{P}(C) &= \int_{\mathcal{X}} \int_{\mathcal{Y}} 1_C(x, y) p(x) q(y) Q_o(dy) P_o(dx) \\ &= \int_{\mathcal{X}} \int_{C_x} q(y) Q_o(dy) p(x) P_o(dx) \\ &= \int_{\mathcal{X}} Q(C_x) p(x) P_o(dx) \\ &= \int_{\mathcal{X}} Q(C_x) P(dx). \end{aligned}$$

More generally, for \mathcal{C} -measurable functions $f : \Omega \rightarrow [0, \infty]$,

$$\begin{aligned} \int_{\Omega} f d\mathbb{P} &= \int_{\Omega} f(x, y) p(x) q(y) \mathbb{P}_o(d(x, y)) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} f(x, y) p(x) q(y) Q_o(dy) P_o(dx) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} f(x, y) Q(dy) p(x) P_o(dx) \\ &= \int_{\mathcal{X}} h(x) P(dx) \end{aligned}$$

with $h(x) := \int_{\mathcal{Y}} f(x, y) Q(dy)$.

The next result concerns uniqueness of stochastic kernels.

Lemma 4.7 (Uniqueness in product measures). *Let P, \tilde{P} be finite measures on $(\mathcal{X}, \mathcal{A})$, and let K, \tilde{K} be stochastic kernels from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ such that*

$$P \otimes K = \tilde{P} \otimes \tilde{K}.$$

Then $P \equiv \tilde{P}$, and for any $B \in \mathcal{B}$, $K(\cdot, B) = \tilde{K}(\cdot, B)$ P -almost everywhere. If \mathcal{B} has a countable generating family \mathcal{E} , then even

$$K(x, \cdot) \equiv \tilde{K}(x, \cdot) \quad \text{for } P\text{-almost all } x \in \mathcal{X}.$$

Proof of Lemma 4.7. The construction of the product measures implies that for any $A \in \mathcal{A}$,

$$P(A) = P \otimes K(A \times \mathcal{Y}) = \tilde{P} \otimes \tilde{K}(A \times \mathcal{Y}) = \tilde{P}(A),$$

i.e. $P \equiv \tilde{P}$. For any fixed $B \in \mathcal{B}$, the functions $K(\cdot, B), \tilde{K}(\cdot, B) \in \mathcal{L}^1(P)$ satisfy

$$\int_A K(x, B) P(dx) = P \otimes K(A \times B) = P \otimes \tilde{K}(A \times B) = \int_A \tilde{K}(x, B) P(dx)$$

for arbitrary $A \in \mathcal{A}$. Hence $K(\cdot, B) = \tilde{K}(\cdot, B)$ P -almost everywhere by Lemma 2.20.

Finally, suppose that \mathcal{E} is a countable family of subsets of \mathcal{Y} generating \mathcal{B} . The family \mathcal{E}' of intersections of finitely many sets in \mathcal{E} is countable as well, see Exercise 4.8. Then the set

$$A_* := \{x \in \mathcal{X} : K(x, B) \neq \tilde{K}(x, B) \text{ for some } B \in \mathcal{E}'\}$$

satisfies $P(A_*) = 0$. But for any fixed $x \in \mathcal{X} \setminus A_*$,

$$\mathcal{D}(x) := \{B \in \mathcal{B} : K(x, B) = \tilde{K}(x, B)\}$$

defines a Dynkin system with $\mathcal{E}' \subset \mathcal{D}(x) \subset \mathcal{B}$. Since \mathcal{E}' is closed under (finite) intersections, the smallest Dynkin system containing \mathcal{E}' is a σ -field, so $\mathcal{D}(x) = \mathcal{B}$ and $K(x, \cdot) \equiv \tilde{K}(x, \cdot)$. \square

Exercise 4.8 (Countable set families). Let \mathcal{E} be a countable family of subsets of \mathcal{Y} .

(a) Show that the family \mathcal{E}' of all intersections of finitely many sets in \mathcal{E} is countable as well, and that \mathcal{E}' is closed under (finite) intersections.

(b) Let \mathcal{B}_o be the family of all sets of the form

$$\bigcup_{i=1}^m \bigcap_{j=1}^n E_{ij}$$

with $m, n \in \mathbb{N}$ and $E_{ij} \in \mathcal{E} \cup \{\mathcal{Y} \setminus E : E \in \mathcal{E}\} \cup \{\emptyset, \mathcal{Y}\}$. Show that \mathcal{B}_o is countable, and that it is the smallest field over \mathcal{Y} containing \mathcal{E} .

4.2 Decomposing Measures on Product Spaces

An interesting question is whether any finite measure \mathbb{P} on (Ω, \mathcal{C}) may be decomposed into the product $P \otimes K$ of a finite measure and a stochastic kernel. Let us start with two examples for such a decomposition.

Example 4.9. Suppose that \mathcal{X} is a countable set and $\mathcal{A} = \mathcal{P}(\mathcal{X})$. Then any finite measure \mathbb{P} on Ω may be represented as $P \otimes K$, where

$$P(\{x\}) := \mathbb{P}(\{x\} \times \mathcal{Y})$$

and

$$K(x, B) := \begin{cases} P(\{x\})^{-1} \mathbb{P}(\{x\} \times B) & \text{if } P(\{x\}) > 0, \\ \mathbb{P}(\mathcal{X})^{-1} \mathbb{P}(\mathcal{X} \times B) & \text{if } P(\{x\}) = 0. \end{cases}$$

Example 4.10. Let L and M be σ -finite measures on $(\mathcal{X}, \mathcal{A})$ and $(\mathcal{Y}, \mathcal{B})$, respectively. Suppose that \mathbb{P} is a finite measure on $(\mathcal{X} \times \mathcal{Y}, \mathcal{A} \otimes \mathcal{B})$ with density $f : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$ with respect to the product measure $L \otimes M$. By Fubini's theorem (Remark 4.6),

$$\mathbb{P}(C) = \int_{\mathcal{X}} \int_{\mathcal{Y}} 1_C(x, y) f(x, y) M(dy) L(dx) = \int_{\mathcal{X}} \int_{C_x} f(x, y) M(dy) L(dx)$$

for arbitrary $C \in \mathcal{C} = \mathcal{A} \otimes \mathcal{B}$. If we define

$$P(A) := \mathbb{P}(A \times \mathcal{Y})$$

for $A \in \mathcal{A}$, then P is a finite measure on $(\mathcal{X}, \mathcal{A})$ with $P(\mathcal{X}) = \mathbb{P}(\Omega)$ and

$$P(A) = \int_A f_1(x) P(dx) \quad \text{with} \quad f_1(x) := \int_{\mathcal{Y}} f(x, y) M(dy).$$

Thus $f_1 = dP/dL$. Note also that $P(\{f_1 = 0\}) = 0 = P(\{f_1 = \infty\})$.

Analogously, $Q(B) := \mathbb{P}(\Omega)^{-1} \mathbb{P}(\mathcal{X} \times B)$ defines a probability measure on $(\mathcal{Y}, \mathcal{B})$ with density $f_2 = dQ/dM$ given by

$$f_2(y) := \mathbb{P}(\Omega)^{-1} \int_{\mathcal{X}} f(x, y) L(dx).$$

For $C \in \mathcal{C}$, we may rewrite $\mathbb{P}(C)$ as

$$\mathbb{P}(C) = \int_{\mathcal{X}} f_1(x) K(x, C_x) L(dx) = \int_{\mathcal{X}} K(x, C_x) P(dx)$$

with

$$K(x, B) := \int_B f_{2|1}(y | x) M(dy)$$

and

$$f_{2|1}(y | x) := \begin{cases} f_1(x)^{-1} f(x, y) & \text{if } 0 < f_1(x) < \infty, \\ f_2(y) & \text{else.} \end{cases}$$

This defines a stochastic kernel K from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ such that $\mathbb{P} = P \otimes K$.

Exercise 4.11 (Decomposition of a probability distribution). Let \mathbb{P} be the uniform distribution on the unit sphere in $\mathbb{R} \times \mathbb{R}$. That is, it describes the distribution of $(\cos(U), \sin(U))$, where $U \sim \text{Unif}[0, 2\pi]$. Determine a decomposition $\mathbb{P} = P \otimes K$ with a probability measure P on \mathbb{R} given by some density, and a stochastic kernel K from \mathbb{R} to \mathbb{R} .

Our final result in this section will show that for “nice” measurable spaces $(\mathcal{Y}, \mathcal{B})$, any finite measure on \mathcal{C} may be decomposed into the product of a finite measure on $(\mathcal{X}, \mathcal{A})$ and a stochastic kernel from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$.

Theorem 4.12 (Existence of stochastic kernels). *Let (\mathcal{Y}, d) be a separable and complete metric space¹, and let $\mathcal{B} = \text{Borel}(\mathcal{Y}, d)$. Then for any finite measure \mathbb{P} on (Ω, \mathcal{C}) there exists a finite measure P on $(\mathcal{X}, \mathcal{A})$ and a stochastic kernel K from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ such that*

$$\mathbb{P} = P \otimes K.$$

¹The corresponding topological space is called a *Polish space*.

For the proof of this theorem, we need some basic facts about finite measures on metric spaces.

Lemma 4.13. *Let (\mathcal{Y}, d) be a metric space, and let Q be a finite measure on $\mathcal{B} = \text{Borel}(\mathcal{Y}, d)$.*

(i) *Any set $B \in \mathcal{B}$ may be approximated in the following sense: For each $\epsilon > 0$ there exist a closed set $A \subset \mathcal{Y}$ and an open set $U \subset \mathcal{Y}$ such that $A \subset B \subset U$ and $Q(U \setminus A) \leq \epsilon$.*

(ii) *If (\mathcal{Y}, d) is separable and complete, the closed sets A in part (i) may even be chosen to be compact.*

Proof of Lemma 4.13. The proof of part (i) is left to the reader as an exercise. One can proceed as follows: Let \mathcal{B}' be the set of all $B \in \mathcal{B}$ which may be approximated by closed and open sets as stated. One can show that \mathcal{B}' is a field over \mathcal{Y} . Then one can show that it is even a σ -field over \mathcal{Y} by showing that $\bigcup_{n \geq 1} B_n \in \mathcal{B}'$ for arbitrary sets $B_1 \subset B_2 \subset B_3 \subset \dots$ in \mathcal{B}' . Any closed set $A \subset \mathcal{Y}$ belongs to \mathcal{B} . To this end consider the neighborhoods $U_\epsilon(A) := \{y \in \mathcal{Y} : d(y, A) < \epsilon\}$ with $d(y, A) := \inf\{d(y, z) : z \in A\}$. Consequently, \mathcal{B}' is a σ -field containing all closed sets and contained in \mathcal{B} . But this implies that $\mathcal{B}' = \mathcal{B}$.

As to part (ii), let $\{y_1, y_2, y_3, \dots\}$ be a dense subset of \mathcal{Y} . For any fixed $\epsilon > 0$ and arbitrary $k \in \mathbb{N}$ let $N(\epsilon, k) \in \mathbb{N}$ such that

$$Q\left(\mathcal{Y} \setminus \bigcup_{n=1}^{N(\epsilon, k)} B(y_n, 1/k)\right) \leq 2^{-k} \epsilon,$$

where $B(y, \delta)$ is the closed ball with center y and radius $\delta \geq 0$. Then

$$K_\epsilon := \bigcap_{k=1}^{\infty} \bigcup_{n=1}^{N(k, \epsilon)} B(y_n, 1/k)$$

is a closed subset of \mathcal{Y} such that $Q(\mathcal{Y} \setminus K_\epsilon) \leq \epsilon$. Indeed, the set K_ϵ is even compact. To show this, let $(x_m)_{m \geq 1}$ be an arbitrary sequence in K_ϵ . Since K_ϵ is closed and (\mathcal{Y}, d) is complete, it suffices to show that $(x_m)_{m \geq 1}$ has a Cauchy subsequence. There exists an index $n(1) \in \{1, \dots, N(\epsilon, k)\}$ such that $M_1 := \{m \geq 1 : x_m \in B(y_{n(1)}, 1)\}$ is infinite. Suppose that for some $k \geq 1$ we have chosen an infinite set $M_k \subset \mathbb{N}$ and an index $n(k) \in \{1, \dots, N(\epsilon, k)\}$ such that $x_m \in B(y_{n(k)}, 1/k)$ for $m \in M_k$. Then for a suitable index $n(k+1) \in \{1, \dots, N(\epsilon, k+1)\}$, the set

$$M_{k+1} := \{m \in M_k : x_m \in B(y_{n(k+1)}, 1/(k+1))\}$$

is infinite too. This leads to infinite subsets $M_1 \supset M_2 \supset M_3 \supset \dots$ of \mathbb{N} such that for each $k \geq 1$,

$$x_m \in B(y_{n(k)}, 1/k) \quad \text{for all } m \in M_k.$$

Now let $m(k)$ be the k -th smallest element of M_k . Then $m(1) < m(2) < m(3) < \dots$, and for integers $1 \leq k \leq \ell$, $m(k), m(\ell) \in M_k$, whence

$$d(x_k, x_\ell) \leq d(x_k, y_{n(k)}) + d(x_\ell, y_{n(k)}) \leq 2/k.$$

This implies that $(x_{m(k)})_{k \geq 1}$ is a Cauchy subsequence of $(x_m)_{m \geq 1}$, whence K_ϵ is compact.

For any set $B \in \mathcal{B}$ and $\epsilon > 0$ there exist a closed set A and an open set U such that $A \subset B \subset U$ and $Q(U \setminus A) \leq \epsilon/2$. But then, $\tilde{A} := A \cap K_{\epsilon/2}$ is a compact subset of B such that $Q(U \setminus \tilde{A}) \leq \epsilon$. \square

Proof of Theorem 4.12. Without loss of generality let \mathbb{P} be a probability measure on (Ω, \mathcal{C}) . For $A \in \mathcal{A}$ and $B \in \mathcal{B}$ let $P(A) := \mathbb{P}(A \times \mathcal{Y})$ and $Q(B) := \mathbb{P}(\mathcal{X} \times B)$. This defines probability measures P and Q on $(\mathcal{X}, \mathcal{A})$ and $(\mathcal{Y}, \mathcal{B})$, respectively.

Since (\mathcal{Y}, d) is separable, there exists a countable dense subset \mathcal{Y}_o of \mathcal{Y} . Let \mathcal{E} be the family of all open balls with center in \mathcal{Y}_o and rational radius. This family \mathcal{E} is countable, and any open subset of \mathcal{Y} is the union of balls in \mathcal{E} . As shown in Exercise 4.8, the smallest field \mathcal{B}_o over \mathcal{Y} containing \mathcal{E} is countable, too. Hence, $\mathcal{B} = \sigma(\mathcal{E}) = \sigma(\mathcal{B}_o)$. Since (\mathcal{Y}, d) is complete, for each $B_o \in \mathcal{B}_o$ there exist compact sets

$$D_1(B_o) \subset D_2(B_o) \subset D_3(B_o) \subset \cdots \subset B_o$$

such that $\lim_{k \rightarrow \infty} Q(D_k(B_o)) = Q(B_o)$, see Lemma 4.13. Finally, we consider the smallest field $\bar{\mathcal{B}}_o$ over \mathcal{Y} containing \mathcal{B}_o and $\{D_k(B_o) : B_o \in \mathcal{B}_o, k \in \mathbb{N}\}$, which is also countable.

For each $B \in \bar{\mathcal{B}}_o$ let $K(\cdot, B)$ be an explicit version of $\mathbb{E}(1_B(Y) | X)$, where $X(x, y) := x$ and $Y(x, y) := y$ for $(x, y) \in \Omega$. That means, for arbitrary $A \in \mathcal{A}$,

$$\mathbb{P}(A \times B) = \mathbb{E}(1_A(X)1_B(Y)) = \mathbb{E}(1_A(X)K(X, B)) = \int_A K(x, B) P(dx).$$

The properties of conditional expectations imply that P -almost everywhere,

$$\begin{aligned} K(\cdot, \emptyset) &= 0, \\ K(\cdot, \mathcal{Y}) &= 1, \\ K(\cdot, B \cup B') &= K(\cdot, B) + K(\cdot, B') \quad \text{for disjoint sets } B, B' \in \bar{\mathcal{B}}_o. \end{aligned}$$

The collection of all these equations is countable, so there exists a set $A_1 \in \mathcal{A}$ with $P(A_1) = 0$ such that for all $x \in \mathcal{X} \setminus A_1$, the mapping $K(x, \cdot)$ defines a probability content on $\bar{\mathcal{B}}_o$, i.e. $0 \leq K(x, \cdot) \leq 1$ with $K(x, \emptyset) = 0$, $K(x, \mathcal{Y}) = 1$ and $K(x, B \cup B') = K(x, B) + K(x, B')$ for disjoint sets $B, B' \in \bar{\mathcal{B}}_o$.

For any set B_o in the smaller set algebra \mathcal{B}_o , we know that

$$\lim_{k \rightarrow \infty} Q(D_k(B_o)) = Q(B_o).$$

But $Q(D_k(B_o)) = \int K(x, D_k(B_o)) P(dx)$, $Q(B_o) = \int K(x, B_o) P(dx)$, and for any $x \in \mathcal{X} \setminus A_1$, the sequence $(K(x, D_k(B_o)))_{k \geq 1}$ is increasing with a limit $h(x, B_o) \leq K(x, B_o)$. By monotone convergence, $\int h(x, B_o) P(dx) = Q(B_o) = \int K(x, B_o) P(dx)$, whence $K(\cdot, B_o) = h(\cdot, B_o)$ for P -almost all $x \in \mathcal{X}$. Since \mathcal{B}_o is countable, there exists a set $A_2 \in \mathcal{A}$ with $A_2 \subset \mathcal{X} \setminus A_1$ and $P(A_2) = 0$ such that for any $x \in \mathcal{X} \setminus (A_1 \cup A_2)$,

$$\lim_{k \rightarrow \infty} K(x, D_k(B_o)) = K(x, B_o) \quad \text{for each } B_o \in \mathcal{B}_o.$$

These properties imply that for each $x \in \mathcal{X} \setminus (A_1 \cup A_2)$, the mapping $K(x, \cdot)$ is a probability measure on the smaller set algebra \mathcal{B}_o . To prove this it suffices to show that for any sequence $(B_n)_n$ of sets $B_1 \supset B_2 \supset B_3 \supset \cdots$ in \mathcal{B}_o with $\bigcap_{n \geq 1} B_n = \emptyset$,

$$\lim_{n \rightarrow \infty} K(x, B_n) = 0.$$

Indeed, for any fixed $\epsilon > 0$ there exist indices $k(n) \geq 1$ such that

$$K(x, B_n \setminus D_n) = K(x, B_n) - K(x, D_n) \leq 2^{-n}\epsilon \quad \text{with } D_n := D_{k(n)}(B_n).$$

Note that D_n is a subset of B_n , so $\bigcap_{n \geq 1} D_n = \emptyset$. Since all sets D_n are compact, there exists an integer $n_o \geq 1$ such that $\bigcap_{m=1}^{n_o} D_m = \emptyset$. Consequently, for $n \geq n_o$,

$$\begin{aligned} K(x, B_n) &\leq K(x, B_{n_o}) = K(x, B_{n_o} \setminus \bigcap_{m=1}^{n_o} D_m) \\ &= K\left(x, \bigcup_{m=1}^{n_o} B_{n_o} \setminus D_m\right) \\ &\leq K\left(x, \bigcup_{m=1}^{n_o} B_m \setminus D_m\right) \\ &\leq \sum_{m=1}^{n_o} K(x, B_m \setminus D_m) \leq \sum_{m=1}^{n_o} 2^{-m}\epsilon < \epsilon. \end{aligned}$$

This shows that $K(x, B_n) \rightarrow 0$ as $n \rightarrow \infty$.

For each $x \in \mathcal{X} \setminus (A_1 \cup A_2)$, there exists a unique extension of $K(x, \cdot)$ to a probability measure $K^*(x, \cdot)$ on $(\mathcal{Y}, \mathcal{B})$. If we set

$$K^*(x, B) := Q(B) \quad \text{for } x \in A_1 \cup A_2 \text{ and } B \in \mathcal{B},$$

then $K^* : \mathcal{X} \times \mathcal{B}$ has the following properties:

- For each $x \in \mathcal{X}$, $K^*(x, \cdot)$ is a probability measure on $(\mathcal{Y}, \mathcal{B})$;
- for each $B_o \in \mathcal{B}_o$, $K^*(\cdot, B_o)$ is \mathcal{A} -measurable and a version of $\mathbb{E}(1_{B_o}(Y) | X)$.

But the set of all $B \in \mathcal{B}$ such that $K^*(\cdot, B)$ is \mathcal{A} -measurable is easily seen to be a Dynkin system. Since it contains the set algebra \mathcal{B}_o , it coincides with \mathcal{B} . Hence K^* is indeed a stochastic kernel, and the distribution $P \otimes K^*$ coincides with \mathbb{P} on the family $\{A \times B_o : A \in \mathcal{A}, B_o \in \mathcal{B}_o\}$. The latter family is closed under intersections and generates \mathcal{C} , so $P \otimes K^* \equiv \mathbb{P}$. \square

4.3 Conditional Expectations and Distributions

So far, we know conditional expectations as functions with certain properties. But the term “expectation” refers to integrals, so an obvious question is whether a conditional expectation may be viewed as an integral with respect to some “conditional distribution”. As shown below, the answer is yes.

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and consider random variables $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \mathcal{A})$ and $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \mathcal{B})$. This gives rise to three different distributions:

$$\begin{aligned} \mathcal{A} \ni A &\mapsto \mathbb{P}^X(A) := \mathbb{P}(X \in A) && \text{(distribution of } X), \\ \mathcal{B} \ni B &\mapsto \mathbb{P}^Y(B) := \mathbb{P}(Y \in B) && \text{(distribution of } Y), \\ \mathcal{A} \otimes \mathcal{B} \ni C &\mapsto \mathbb{P}^{(X,Y)}(C) := \mathbb{P}((X, Y) \in C) && \text{(joint distribution of } X \text{ and } Y). \end{aligned}$$

Now suppose that there exists a stochastic kernel K from $(\mathcal{X}, \mathcal{A})$ to $(\mathcal{Y}, \mathcal{B})$ such that

$$\mathbb{P}^{(X,Y)} = \mathbb{P}^X \otimes K.$$

Then $K(x, \cdot)$ may be interpreted as the *conditional distribution of Y given $X = x$* . Indeed, let $f : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty]$ be $\mathcal{A} \otimes \mathcal{B}$ -measurable. Then it follows from Theorem 4.5 that

$$\mathbb{E} f(X, Y) = \int_{\mathcal{X} \times \mathcal{Y}} f d\mathbb{P}^{(X,Y)} = \int_{\mathcal{X}} h d\mathbb{P}^X$$

with the \mathcal{A} -measurable function $h : \mathcal{X} \rightarrow [0, \infty]$,

$$h(x) := \int_{\mathcal{Y}} f(x, y) K(x, dy).$$

The same conclusion is true if $f \in \mathcal{L}^1(\mathbb{P}^{(X,Y)})$, the only caveat being that the integral $h(x)$ may not exist for points x in a set $A(f) \in \mathcal{A}$ such that $\mathbb{P}^X(A(f)) = 0$.

In particular, for $g \in \mathcal{L}^1(\mathbb{P}^Y)$ and arbitrary sets $A \in \mathcal{A}$,

$$\begin{aligned} \int_{X^{-1}(A)} g(Y) d\mathbb{P} &= \int_{\Omega} 1_A(X) g(Y) d\mathbb{P} \\ &= \int_{\mathcal{X} \times \mathcal{Y}} 1_A(x) g(y) \mathbb{P}^{(X,Y)}(d(x, y)) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} 1_A(x) g(y) K(x, dy) \mathbb{P}^X(dx) \\ &= \int_A \int_{\mathcal{Y}} g(y) K(x, dy) \mathbb{P}^X(dx), \end{aligned}$$

so

$$x \mapsto \int_{\mathcal{Y}} g(y) K(x, dy)$$

is a version of $\mathbb{E}(g(Y) | X)$. One often writes

$$\left. \begin{array}{l} \mathbb{P}(Y \in B | X = x) \\ \mathbb{P}^Y(B | X = x) \end{array} \right\} \text{ instead of } K(x, B),$$

$$\mathbb{E}(g(Y) | X = x) \text{ instead of } \int_{\mathcal{Y}} g(y) K(x, dy).$$

Optimal predictions. In the special case that (\mathcal{Y}, d) is a separable, complete metric space and $\mathcal{B} = \text{Borel}(\mathcal{Y}, d)$, it follows from Theorem 4.12 that a kernel K with the stated properties does exist. In particular, let $\mathcal{Y} = \mathbb{R}$. If $Y \in \mathcal{L}^1(\mathbb{P})$, then a version of $\mathbb{E}(Y | X)$ is given by

$$\mathbb{E}(Y | X = x) := \int_{\mathbb{R}} y K(x, dy).$$

Note that

$$\infty > \int |Y| d\mathbb{P} = \int_{\mathcal{X}} |y| K(x, dy) \mathbb{P}^X(dx),$$

so we may redefine $K(x, \cdot) := \mathbb{P}^Y$ for the exceptional points $x \in \mathcal{X}$ such that $\int |y| K(x, dy) = \infty$. With the kernel K at hand, we can solve various prediction problems. The goal is to find a *predictor* $g(X)$ of Y , determined by a measurable function $g : (\mathcal{X}, \mathcal{A}) \rightarrow \mathbb{R}$, such that the approximation error $Y - g(X)$ is “small”. The solution of this prediction problem depends on the way how we measure the approximation error precisely.

Case 1: Mean squared prediction error. Suppose that $Y \in \mathcal{L}^2(\mathbb{P})$. If we want to minimize $\mathbb{E}((Y - g(X))^2)$, one can deduce from the considerations in Section 3.2 that $\mathbb{E}(Y | X)$ solves this problem. But let us follow a more direct route using the kernel K . Recall first that for any constant $q \in \mathbb{R}$, the expectation of $(Y - q)^2$ equals $\text{Var}(Y) + (\mathbb{E}(Y) - q)^2$, where $\text{Var}(Y) = \mathbb{E}((Y - \mathbb{E}(Y))^2)$. Note also that

$$\infty > \int Y^2 d\mathbb{P} = \int_{\mathcal{X}} \int_{\mathbb{R}} y^2 K(x, dy) \mathbb{P}^X(dx),$$

so we may redefine $K(x, \cdot) := \mathbb{P}^Y$ whenever $\int y^2 K(x, dy) = \infty$. Now we can write

$$\begin{aligned} \mathbb{E}((Y - g(X))^2) &= \int_{\mathcal{X}} \int_{\mathbb{R}} (y - g(x))^2 K(x, dy) \mathbb{P}^X(dx) \\ &= \int_{\mathcal{X}} \left(\text{Var}(Y | X = x) + (\mathbb{E}(Y | X = x) - g(x))^2 \right) \mathbb{P}^X(dx), \end{aligned}$$

where

$$\text{Var}(Y | X = x) := \int_{\mathbb{R}} (y - \mathbb{E}(Y | X = x))^2 K(x, dy).$$

Consequently, a predictor $g(X)$ of Y minimizes the mean squared prediction error if and only if $g(X) = \mathbb{E}(Y | X)$ almost surely.

Case 2: Mean absolute prediction error. Suppose we want to minimize $\mathbb{E}|Y - g(X)|$. It is well-known that for a real constant q , the expectation of $|Y - q|$ is minimal if and only if q is a median of (the distribution of) Y , that means $\mathbb{P}(Y < q) \leq 1/2 \leq \mathbb{P}(Y \geq q)$. This leads to an optimal predictor $g(X)$:

$$\mathbb{E}|Y - g(X)| = \int_{\mathcal{X}} \int_{\mathbb{R}} |y - g(x)| K(x, dy) \mathbb{P}^X(dx)$$

is minimal if and only if $g(x)$ is a median for $\mathbb{P}(Y \in \cdot | X = x)$, that means,

$$\mathbb{P}(Y < g(x) | X = x) \leq 1/2 \leq \mathbb{P}(Y \leq g(x) | X = x),$$

for \mathbb{P}^X -almost all $x \in \mathcal{X}$.

Chapter 5

Haar Measure

5.1 Locally Compact Topological Groups

Groups. Recall the definition of a group (\mathcal{X}, \cdot) . This is a set \mathcal{X} with a binary operation

$$\mathcal{X} \times \mathcal{X} \ni (x, y) \mapsto x \cdot y \in \mathcal{X}$$

satisfying the following three conditions:

- (Associativity) For arbitrary $x, y, z \in \mathcal{X}$,

$$x \cdot (y \cdot z) = (x \cdot y) \cdot z.$$

- (Identity element) There exists an element $e \in \mathcal{X}$ such that for arbitrary $x \in \mathcal{X}$,

$$x \cdot e = e \cdot x = x.$$

- (Inverse elements) For each $x \in \mathcal{X}$ there exists an element $x^{-1} \in \mathcal{X}$ such that

$$x \cdot x^{-1} = x^{-1} \cdot x = e.$$

The identity element e is unique, and for each $x \in \mathcal{X}$, its inverse element x^{-1} is unique as well. Instead of $x \cdot y$ one often writes xy .

Exercise 5.1. Let (\mathcal{X}, \cdot) be a group.

(a) For $x, y \in \mathcal{X}$ let $x \star y := yx$. Show that (\mathcal{X}, \star) is also a group with the same identity element and the same inverse elements.

(b) Show that for any fixed $x_o \in \mathcal{X}$, the mapping $x \mapsto x_o x$ as well as the mapping $x \mapsto x x_o$ is bijective from \mathcal{X} to \mathcal{X} .

Locally compact topological groups. Suppose that there is also a metric d on \mathcal{X} such that the following two conditions are satisfied:¹

¹For simplicity of exposition, we restrict our attention to metric spaces, although the subsequent definitions and main results may be extended to non-metric topological spaces.

- The mapping $x \mapsto x^{-1}$, from \mathcal{X} to \mathcal{X} is continuous.
- The mapping $(x, y) \mapsto xy$, from $\mathcal{X} \times \mathcal{X}$ to \mathcal{X} is continuous.²

Then one calls (\mathcal{X}, \cdot, d) a *topological group*.

Suppose that in addition,

- for each $x \in \mathcal{X}$ there exists an $\epsilon > 0$ such that the closed ball $B_\epsilon(x)$ with center x and radius ϵ is compact.

Then one calls (\mathcal{X}, \cdot, d) a *locally compact topological group*. If the group operation and metric are clear from the context, we just talk about the (locally compact) topological group \mathcal{X} .

Now we provide three examples of locally compact topological groups.

Example 5.2 (Euclidean spaces). Let $\mathcal{X} = \mathbb{R}^d$. With the usual addition of vectors and the usual Euclidean distance, \mathbb{R}^d becomes a locally compact topological group which is also commutative.

The following two examples involve sets of invertible matrices in $\mathbb{R}^{d \times d}$ for some $d \in \mathbb{N}$, and the binary operation is matrix multiplication. To stay coherent with the general theory in this section, we denote matrices temporarily with lower-case, boldface letters, and the identity matrix is denoted with e . Note that $\mathbb{R}^{d \times d}$ may be identified with \mathbb{R}^{d^2} by identifying a matrix $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$ having columns $\mathbf{x}_j = (x_{ij})_{i=1}^d \in \mathbb{R}^d$ with the vector $(\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_d^\top)^\top$. Thus the usual Euclidean norm on \mathbb{R}^{d^2} corresponds to Frobenius norm on $\mathbb{R}^{d \times d}$, i.e.

$$\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i,j=1}^d x_{ij}y_{ij} \quad \text{and} \quad \|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

Note also that any other norm on $\mathbb{R}^{d \times d}$ or \mathbb{R}^{d^2} induces the same topology. Finally, a basic fact is that matrix multiplication is a continuous mapping from $\mathbb{R}^{d \times d} \times \mathbb{R}^{d \times d}$ to $\mathbb{R}^{d \times d}$.

Example 5.3 (Linear groups). Let \mathcal{X} be the set of all invertible matrices $\mathbf{x} \in \mathbb{R}^{d \times d}$. With the usual matrix multiplication, \mathcal{X} is a non-commutative³ group. Moreover, the set \mathcal{X} is an open subset of $\mathbb{R}^{d \times d}$, and the mapping $x \mapsto x^{-1}$ is continuous from \mathcal{X} to \mathcal{X} . Indeed, for $\mathbf{x} \in \mathcal{X}$ and $\boldsymbol{\delta} \in \mathbb{R}^{d \times d}$, it is well-known that $\mathbf{x} + \boldsymbol{\delta}$ is invertible, provided that $\|\boldsymbol{\delta}\|$ is sufficiently small, and

$$(\mathbf{x} + \boldsymbol{\delta})^{-1} = \mathbf{x}^{-1} - \mathbf{x}^{-1}\boldsymbol{\delta}\mathbf{x}^{-1} + O(\|\boldsymbol{\delta}\|^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

This follows from $\mathbf{x} + \boldsymbol{\delta} = \mathbf{x}(e + \mathbf{x}^{-1}\boldsymbol{\delta})$ and from von Neumann's series expansion

$$(e + \mathbf{a})^{-1} = \sum_{k=0}^{\infty} (-1)^k \mathbf{a}^k$$

for $\mathbf{a} \in \mathbb{R}^{d \times d}$ with $\|\mathbf{a}\|$ sufficiently small. Hence \mathcal{X} is a topological group.

Example 5.4 (Orthogonal groups). Let $\mathcal{X}_{\text{orth}}$ be the set of orthogonal matrices $\mathbf{x} \in \mathbb{R}^{d \times d}$, that means, $\mathbf{x}^\top \mathbf{x} = e$. This is a subgroup of the set \mathcal{X} of all invertible matrices and thus a locally compact topological group. Indeed, $\mathcal{X}_{\text{orth}}$ is even a compact set and a smooth $d(d-1)/2$ -dimensional

²Here $\mathcal{X} \times \mathcal{X}$ is equipped with the metric $d((x, x'), (y, y')) := \max\{d(x, y), d(x', y')\}$.

³unless $d = 1$

manifold. To see this, note first that

$$\mathcal{X}_{\text{orth}} = F^{-1}(e)$$

with the continuous mapping

$$F : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}, \quad F(\mathbf{x}) := \mathbf{x}^\top \mathbf{x}.$$

Hence, $\mathcal{X}_{\text{orth}}$ is a closed subset of $\mathbb{R}^{d \times d}$. Moreover, with the Frobenius norm $\|\cdot\|$, all matrices $\mathbf{x} \in \mathcal{X}_{\text{orth}}$ satisfy $\|\mathbf{x}\| = \sqrt{d}$, so $\mathcal{X}_{\text{orth}}$ is a closed and bounded subset of $\mathbb{R}^{d \times d}$. Hence it is compact.

That $\mathcal{X}_{\text{orth}}$ is a smooth $d(d-1)$ -dimensional manifold can be verified as follows: Note that

$$\mathbb{R}^{d \times d} = \mathbb{R}_{\text{sym}}^{d \times d} + \mathbb{R}_{\text{skew}}^{d \times d},$$

where $\mathbb{R}_{\text{sym}}^{d \times d}$ and $\mathbb{R}_{\text{skew}}^{d \times d}$ are the linear spaces of symmetric and skew-symmetric matrices in $\mathbb{R}^{d \times d}$, respectively:

$$\begin{aligned} \mathbb{R}_{\text{sym}}^{d \times d} &:= \{\mathbf{a} \in \mathbb{R}^{d \times d} : \mathbf{a}^\top = \mathbf{a}\} && \text{with dimension } d(d+1)/2, \\ \mathbb{R}_{\text{skew}}^{d \times d} &:= \{\mathbf{a} \in \mathbb{R}^{d \times d} : \mathbf{a}^\top = -\mathbf{a}\} && \text{with dimension } d(d-1)/2. \end{aligned}$$

Moreover, $\mathbb{R}_{\text{sym}}^{d \times d} \perp \mathbb{R}_{\text{skew}}^{d \times d}$, and any matrix \mathbf{a} may be written as $\mathbf{a} = \mathbf{a}_{\text{sym}} + \mathbf{a}_{\text{skew}}$ with

$$\mathbf{a}_{\text{sym}} := 2^{-1}(\mathbf{a} + \mathbf{a}^\top) \in \mathbb{R}_{\text{sym}}^{d \times d}, \quad \mathbf{a}_{\text{skew}} := 2^{-1}(\mathbf{a} - \mathbf{a}^\top) \in \mathbb{R}_{\text{skew}}^{d \times d}.$$

Note that $F(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$ defines a continuously differentiable mapping

$$F : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}.$$

For any fixed $\mathbf{x} \in \mathcal{X}_{\text{orth}}$ and arbitrary $\boldsymbol{\delta} \in \mathbb{R}^{d \times d}$,

$$\begin{aligned} F(\mathbf{x} + \boldsymbol{\delta}) &= e + \mathbf{x}^\top \boldsymbol{\delta} + \boldsymbol{\delta}^\top \mathbf{x} + \boldsymbol{\delta}^\top \boldsymbol{\delta} \\ &= e + 2(\mathbf{x}^\top \boldsymbol{\delta})_{\text{sym}} + O(\|\boldsymbol{\delta}\|^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}. \end{aligned}$$

Hence the derivative of F at $\mathbf{x} \in \mathcal{X}_{\text{orth}}$ is given by the linear mapping

$$DF(\mathbf{x}) : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}, \quad DF(\mathbf{x})\mathbf{a} = 2(\mathbf{x}^\top \mathbf{a})_{\text{sym}},$$

and the null space of this mapping equals

$$\{\mathbf{a} \in \mathbb{R}^{d \times d} : DF(\mathbf{x})\mathbf{a} = \mathbf{0}\} = \{\mathbf{x}\mathbf{b} : \mathbf{b} \in \mathbb{R}^{d \times d}, \mathbf{b}_{\text{sym}} = \mathbf{0}\} = \{\mathbf{x}\mathbf{b} : \mathbf{b} \in \mathbb{R}_{\text{skew}}^{d \times d}\}$$

whereas

$$DF(\mathbf{x})\mathbf{x}\mathbf{b} = 2\mathbf{b} \quad \text{for all } \mathbf{b} \in \mathbb{R}_{\text{sym}}^{d \times d}.$$

Consequently, the linear mapping $DF(\mathbf{x})$ has full rank $d(d+1)/2$, and the implicit function theorem shows that $\mathcal{X}_{\text{orth}}$ is a differentiable manifold of dimension $d(d-1)/2$. Its tangent space at any point $\mathbf{x} \in \mathcal{X}_{\text{orth}}$ is the null space of $DF(\mathbf{x})$. That means, for small $\epsilon > 0$,

$$\{\mathbf{y} \in \mathcal{X}_{\text{orth}} : \|\mathbf{y} - \mathbf{x}\| < \epsilon\} \approx \{\mathbf{x} + \mathbf{x}\mathbf{b} : \mathbf{b} \in \mathbb{R}_{\text{skew}}^{d \times d}, \|\mathbf{b}\| < \epsilon\}.$$

Exercise 5.5 (An embedding for $\mathcal{X}_{\text{orth}}$). Consider the mapping $\Psi : \mathcal{X} \rightarrow \mathbb{R}^{d \times d}$,

$$\Psi(\mathbf{x}) := \mathbf{x}(\mathbf{x}^\top \mathbf{x})^{-1/2}.$$

Show that Ψ has the following five properties:

- (a) $\Psi(\mathbf{x}) \in \mathcal{X}_{\text{orth}}$ for all $\mathbf{x} \in \mathcal{X}$;
- (b) $\Psi(\mathbf{x}) = \mathbf{x}$ if and only if $\mathbf{x} \in \mathcal{X}_{\text{orth}}$;
- (c) $\Psi(\mathbf{x}\mathbf{y}) = \mathbf{x}\Psi(\mathbf{y})$ for all $\mathbf{x} \in \mathcal{X}_{\text{orth}}$ and $\mathbf{y} \in \mathcal{X}$;
- (d) $\Psi(\mathbf{e} + \boldsymbol{\delta}) = \mathbf{e}$ for all $\boldsymbol{\delta} \in \mathbb{R}_{\text{sym}}^{d \times d}$ such that $\lambda_{\min}(\boldsymbol{\delta}) > -1$;
- (e) for $\boldsymbol{\delta} \in \mathbb{R}^{d \times d}$ such that $\mathbf{e} + \boldsymbol{\delta} \in \mathcal{X}$,

$$\Psi(\mathbf{e} + \boldsymbol{\delta}) = \mathbf{e} + \boldsymbol{\delta}_{\text{skew}} + O(\|\boldsymbol{\delta}\|^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

In part (e) one should use (and possibly verify) the fact that for $\mathbf{a} \in \mathbb{R}_{\text{sym}}^{d \times d}$ with minimal eigenvalue $\lambda_{\min}(\mathbf{a}) > -1$ and any fixed $b \in \mathbb{R}$,

$$(\mathbf{e} + \mathbf{a})^b = \mathbf{e} + b\mathbf{a} + O(\|\mathbf{a}\|^2) \quad \text{as } \mathbf{a} \rightarrow \mathbf{0}.$$

5.2 Left- and Right-Invariant Measures

Throughout this section let \mathcal{X} be a locally compact topological group. In what follows, for $A, B \subset \mathcal{X}$ we write

$$AB := \{ab : a \in A, b \in B\} \quad \text{and} \quad B^{-1} := \{b^{-1} : b \in B\}.$$

Moreover, for $x \in \mathcal{X}$ we set $xB := \{xb : b \in B\}$ and $Ax := \{ax : a \in A\}$. We also write

$$\begin{aligned} \mathcal{C} &:= \{\text{compact subsets of } \mathcal{X}\}, \\ \mathcal{U} &:= \{\text{open subsets of } \mathcal{X}\}, \\ \mathcal{B} &:= \{\text{Borel subsets of } \mathcal{X}\}. \end{aligned}$$

Our first main result is about the existence and (essential) uniqueness of a left-invariant measure on (the Borel subsets of) \mathcal{X} .

Theorem 5.6 (Haar–Weil). *There exists a measure μ on $(\mathcal{X}, \mathcal{B})$ with the following five properties:*

$$(5.1) \quad \mu(xB) = \mu(B) \quad \text{for all } x \in \mathcal{X}, B \in \mathcal{B},$$

$$(5.2) \quad \mu(B) = \inf_{U \in \mathcal{U}: U \supset B} \mu(U) \quad \text{for all } B \in \mathcal{B},$$

$$(5.3) \quad \mu(U) = \sup_{C \in \mathcal{C}: C \subset U} \mu(C) \quad \text{for all } U \in \mathcal{U},$$

$$(5.4) \quad \mu(C) < \infty \quad \text{for all } C \in \mathcal{C},$$

$$(5.5) \quad \mu(U) > 0 \quad \text{for all } U \in \mathcal{U} \setminus \{\emptyset\}.$$

If $\tilde{\mu}$ is another measure with these properties, then there exists a constant $\gamma > 0$ such that $\tilde{\mu} = \gamma\mu$.

Definition 5.7 (Left Haar measure). A measure μ with the properties listed in Theorem 5.6 is called a *left Haar measure on \mathcal{X}* .

The proof presented here is essentially the one of Haar (1933). Since it uses Tikhonov's theorem and thus the axiom of choice, some mathematicians came up with constructive proofs, see Alfsen (1963) and the references therein.

Proof of Theorem 5.6. In this proof, some facts are stated without proof, and the reader should fill in these details. For arbitrary sets $A, B \subset \mathcal{X}$ with $B \neq \emptyset$ we define a “covering number”

$$N(A|B) := \min\{\#\mathcal{X}_o : \mathcal{X}_o \subset \mathcal{X}, A \subset \mathcal{X}_o B\} \in \mathbb{N}_0 \cup \{\infty\}.$$

These covering numbers have the following properties:

$$N(B|B) = 1.$$

$$N(yA|zB) = N(A|B) \text{ for arbitrary } y, z \in \mathcal{X}.$$

$$N(A|B) > 0 \text{ if and only if } A \neq \emptyset.$$

$$N(A|B) < \infty \text{ if } A \in \mathcal{C} \text{ and } B \text{ has non-empty interior.}$$

Furthermore, for $A_1, A_2 \subset \mathcal{X}$,

$$\begin{aligned} N(A_1|B) &\leq N(A_2|B) \quad \text{if } A_1 \subset A_2, \\ N(A_1 \cup A_2|B) &\begin{cases} \leq N(A_1|B) + N(A_2|B), \\ = N(A_1|B) + N(A_2|B) \end{cases} \quad \text{if } A_1 B^{-1} \cap A_2 B^{-1} = \emptyset. \end{aligned}$$

Finally, if B' is another non-empty subset of \mathcal{X} , then

$$(5.6) \quad N(A|B') \leq N(A|B)N(B|B').$$

The rough idea behind Haar's and Weil's construction is that $N(\cdot|B)$ behaves almost like a left-invariant measure if B is a “small” neighborhood of the identity element e . To make this precise, let B_ϵ be the closed ball around e with radius $\epsilon > 0$,

$$B_\epsilon := \{x \in \mathcal{X} : d(x, e) \leq \epsilon\}.$$

Without loss of generality let B_1 be compact⁴. For $0 < \epsilon \leq 1$ we define $\lambda_\epsilon : \mathcal{C} \rightarrow [0, \infty)$ via

$$\lambda_\epsilon(C) := \frac{N(C|B_\epsilon)}{N(B_1|B_\epsilon)}.$$

It follows from the properties of $N(\cdot|\cdot)$ that $\lambda_\epsilon(C)$ is well-defined in $[0, \infty)$ for any $C \in \mathcal{C}$, and

$$\lambda_\epsilon(xC) = \lambda_\epsilon(C) \quad \text{for arbitrary } x \in \mathcal{X}, C \in \mathcal{C}.$$

Moreover, (5.6) implies that

$$\lambda_\epsilon(C) \begin{cases} \leq \frac{N(C|B_1)N(B_1|B_\epsilon)}{N(B_1|B_\epsilon)} = N(C|B_1) \\ \geq \frac{N(C|B_\epsilon)}{N(B_1|A)N(C|B_\epsilon)} = N(B_1|A)^{-1} \end{cases}$$

⁴Otherwise, replace $d(\cdot, \cdot)$ with $\epsilon_o^{-1}d(\cdot, \cdot)$ for sufficiently small $\epsilon_o > 0$.

with the conventions $N(B_1 | \emptyset) := \infty$ and $1/\infty := 0$. In particular, $\lambda_\epsilon(B_1) = 1$. Finally, for $C_1, C_2 \in \mathcal{C}$,

$$(5.7) \quad \begin{aligned} \lambda_\epsilon(C_1) &\leq \lambda_\epsilon(C_2) \quad \text{if } C_1 \subset C_2, \\ \lambda_\epsilon(C_1 \cup C_2) &\begin{cases} \leq \lambda_\epsilon(C_1) + \lambda_\epsilon(C_2), \\ = \lambda_\epsilon(C_1) + \lambda_\epsilon(C_2) \end{cases} \quad \text{if } C_1 B_\epsilon^{-1} \cap C_2 B_\epsilon^{-1} = \emptyset, \end{aligned}$$

Note that each λ_ϵ , viewed as a tuple $(\lambda_\epsilon(C))_{C \in \mathcal{C}}$, is a point in the infinite Cartesian product

$$\prod_{C \in \mathcal{X}} K_C$$

with the compact intervals $K_C := [N(B_1 | C)^{-1}, N(C | B_1)]$. It follows from Tikhonov's theorem that this product, equipped with the corresponding product topology, is a compact topological space. Hence, the functions λ_ϵ , $\epsilon \in (0, 1]$, have a cluster point $\lambda \in \prod_{C \in \mathcal{C}} K_C$ as $\epsilon \downarrow 0$. Interpreting λ as a function $\lambda : \mathcal{C} \rightarrow [0, \infty)$, this cluster point property means that for any finite set $\mathcal{C}_o \subset \mathcal{C}$,

$$\liminf_{\epsilon \downarrow 0} \max_{C \in \mathcal{C}_o} |\lambda_\epsilon(C) - \lambda(C)| = 0.$$

In particular, λ inherits various nice properties of the functions λ_ϵ :

$$\begin{aligned} \lambda(C) &\in K_C \quad \text{for any } C \in \mathcal{C}, \\ \lambda(B_1) &= 1, \\ \lambda(xC) &= \lambda(C) \quad \text{for arbitrary } x \in \mathcal{X}, C \in \mathcal{C}, \end{aligned}$$

and for $C_1, C_2 \in \mathcal{C}$,

$$\begin{aligned} \lambda(C_1) &\leq \lambda(C_2) \quad \text{if } C_1 \subset C_2, \\ \lambda(C_1 \cup C_2) &\begin{cases} \leq \lambda(C_1) + \lambda(C_2), \\ = \lambda(C_1) + \lambda(C_2) \end{cases} \quad \text{if } C_1 \cap C_2 = \emptyset. \end{aligned}$$

The latter equality follows from (5.7) and the fact that for disjoint compact sets C_1 and C_2 , the neighborhoods $C_1 B_\epsilon^{-1}$ and $C_2 B_\epsilon^{-1}$ are disjoint for sufficiently small $\epsilon > 0$.

Now we are ready to define the measure μ . At first we set

$$\mu(U) := \sup_{C \in \mathcal{C}: C \subset U} \lambda(C) \quad \text{for } U \in \mathcal{U}.$$

This definition implies that

$$\mu(\emptyset) = 0 \quad \text{and} \quad \mu(U) \leq \mu(U') \quad \text{for open sets } U \subset U' \subset \mathcal{X}.$$

This allows us to extend μ to a function on $\mathcal{P}(\mathcal{X})$ via

$$\mu(M) := \inf_{U \in \mathcal{U}: U \supset M} \mu(U) \quad \text{for } M \subset \mathcal{X}.$$

Note that

$$\mu(M) \leq \mu(M') \quad \text{for } M \subset M' \subset \mathcal{X}.$$

Since $\lambda(\cdot)$ is left-invariant, the same is true for μ , that means,

$$\mu(xM) = \mu(M) \quad \text{for all } x \in \mathcal{X} \text{ and } M \subset \mathcal{X}.$$

Note also that

$$\lambda(C) \leq \inf_{U \in \mathcal{U}: C \subset U} \mu(U) = \mu(C) \quad \text{for all } C \in \mathcal{C}.$$

The first step follows from the definition of $\mu(U)$ for open sets U , the second step is just the definition of $\mu(C)$. Together with the monotonicity of $\mu(\cdot)$ this implies that

$$\mu(U) = \sup_{C \in \mathcal{C}: C \subset U} \mu(C) \quad \text{for all } U \in \mathcal{U}.$$

To show that $\mu > 0$ on $\mathcal{U} \setminus \{\emptyset\}$ and $\mu < \infty$ on \mathcal{C} , consider arbitrary sets $C \in \mathcal{C} \setminus \{\emptyset\}$ and $U \in \mathcal{U}$ such that $C \subset U$. For each $x \in C$ there exist $U_x \in \mathcal{U}$ and $B_x \in \mathcal{C}$ such that $x \in U_x \subset B_x \subset U$. Then $(U_x)_{x \in C}$ defines a covering of C with open sets, so by compactness of C , there exists a finite set $C_o \subset C$ such that

$$C \subset \tilde{U} := \bigcup_{x \in C_o} U_x \subset \tilde{C} := \bigcup_{x \in C_o} B_x \subset U.$$

Note that $\tilde{U} \in \mathcal{U}$ and $\tilde{C} \in \mathcal{C}$. Hence,

$$\mu(C) \leq \mu(\tilde{U}) \leq \lambda(\tilde{C}) < \infty$$

and

$$\mu(U) \geq \lambda(\tilde{C}) \geq N(B_1 | \tilde{C})^{-1} > 0,$$

because \tilde{C} has non-empty interior.

It remains to show that μ defines a measure on \mathcal{B} . To this end, we apply Carathéodory's theory of outer measures. It suffices to show that μ defines an outer measure on \mathcal{X} and that any open set $U \in \mathcal{U}$ is μ -measurable in the sense that

$$(5.8) \quad \mu(M) \geq \mu(M \cap U) + \mu(M \setminus U) \quad \text{for any } M \subset \mathcal{X}.$$

Proof that μ is an outer measure: Let $(M_n)_{n \geq 1}$ be an arbitrary sequence of sets $M_n \subset \mathcal{X}$, and let $M \subset \bigcup_{n \geq 1} M_n$. We have to show that $\mu(M) \leq \sum_{n \geq 1} \mu(M_n)$. We may assume that the right hand side is finite, because otherwise the claim is trivial. For arbitrary fixed $\epsilon > 0$ let $U_n \in \mathcal{U}$ with $M_n \subset U_n$ and $\mu(U_n) \leq \mu(M_n) + 2^{-n}\epsilon$. Then $U := \bigcup_{n \geq 1} U_n$ is an open set containing M , whence

$$\mu(M) \leq \mu(U) \quad \text{while} \quad \sum_{n \geq 1} \mu(U_n) \leq \sum_{n \geq 1} \mu(M_n) + \epsilon.$$

Consequently, it suffices to show that $\mu(U) \leq \sum_{n \geq 1} \mu(U_n)$. For any compact set $C \subset U$ and $x \in C$ there exist $n(x) \in \mathbb{N}$ and a compact neighborhood B_x of x such that $B_x \subset U_{n(x)}$. But compactness of C implies that for some finite set $C_o \subset C$,

$$C \subset \bigcup_{x \in C_o} B_x = \bigcup_{n=1}^m \bigcup_{x \in C_o: n(x)=n} B_x = \bigcup_{n=1}^m C_n$$

with $m := \max\{n(x) : x \in C_o\}$ and the compact sets

$$C_n := \bigcup_{x \in C_o : n(x)=n} B_x \subset U_n.$$

Consequently, the properties of λ imply that

$$\lambda(C) \leq \lambda\left(\bigcup_{n=1}^m C_n\right) \leq \sum_{n=1}^m \lambda(C_n) \leq \sum_{n=1}^m \mu(U_n),$$

and for $\lambda(C) \uparrow \mu(U)$ the asserted inequality follows.

Proof of (5.8): It suffices to consider the case $\mu(M) < \infty$. For any fixed $\epsilon > 0$, there exists an open set $V \supset M$ such that $\mu(V) \leq \mu(M) + \epsilon$. On the other hand,

$$\mu(M \cap U) \leq \mu(V \cap U) \leq \lambda(C) + \epsilon$$

for some compact set $C \subset V \cap U$. But then

$$\mu(M \setminus U) \leq \mu(V \setminus U) = \mu(V \setminus (V \cap U)) \leq \mu(V \setminus C) \leq \lambda(\tilde{C}) + \epsilon$$

for some compact set $\tilde{C} \subset V \setminus C$. Hence,

$$\mu(M) + \epsilon \geq \mu(V) \geq \lambda(C \cup \tilde{C}) = \lambda(C) + \lambda(\tilde{C}) \geq \mu(M \cap U) + \mu(M \setminus U) - 2\epsilon.$$

As $\epsilon \downarrow 0$, this yields (5.8).

It remains to prove uniqueness of μ up to positive multiples. To this end, let $\tilde{\mu}$ be a second measure on \mathcal{B} with the stated properties. It suffices to show that $\tilde{\mu} = \gamma\mu$ on \mathcal{C}_* for some constant $\gamma > 0$, where \mathcal{C}_* is the set of compact subsets of \mathcal{X} with non-empty interior. To this end, consider arbitrary compact sets $A, B \subset \mathcal{X}$. Then

$$\begin{aligned} \mu(A)\tilde{\mu}(B) &= \int \mu(A)1_B(y) \tilde{\mu}(dy) \\ &= \int \mu(yA)1_B(y) \tilde{\mu}(dy) && \text{(left-invariance of } \mu) \\ &= \int \int 1_{yA}(x)1_B(y) \mu(dx) \tilde{\mu}(dy) \\ &= \int \int 1_{yA}(x)1_B(y) \tilde{\mu}(dy) \mu(dx) && \text{(Fubini's theorem)} \\ &= \int \int 1_{xA^{-1}}(y)1_B(y) \tilde{\mu}(dy) \mu(dx) && (x \in yA \text{ iff } y \in xA^{-1}) \\ &= \int \tilde{\mu}((xA^{-1}) \cap B) \mu(dx) \\ &= \int \tilde{\mu}(A^{-1} \cap (x^{-1}B)) \mu(dx) && \text{(left-invariance of } \tilde{\mu}) \\ &\leq \int 1_{BA}(x)\tilde{\mu}(A^{-1}) \mu(dx) && (A^{-1} \cap (x^{-1}B) = \emptyset \text{ if } x \notin BA) \\ &= \mu(BA)\tilde{\mu}(A^{-1}). \end{aligned}$$

Note that the application of Fubini's theorem was justified because the integrand, $1_{yA}(x)1_B(y)$, equals zero if $x \notin BA$ or $y \notin B$, so μ and $\tilde{\mu}$ may be viewed temporarily as finite measures on the compact sets BA and B , respectively. Interchanging the roles of μ and $\tilde{\mu}$ yields the inequalities

$$\mu(A)\tilde{\mu}(B) \leq \mu(BA)\tilde{\mu}(A^{-1}) \quad \text{and} \quad \tilde{\mu}(A)\mu(C) \leq \tilde{\mu}(CA)\mu(A^{-1})$$

for $A, B, C \in \mathcal{C}$. Specifically let B, C be arbitrary fixed sets in \mathcal{C}_* , and let $A := B_\epsilon \cap B_\epsilon^{-1}$ for some $\epsilon \in (0, 1]$, so $A^{-1} = A$. Then

$$\frac{\tilde{\mu}(B)}{\mu(BA)} \leq \frac{\tilde{\mu}(A)}{\mu(A)} \leq \frac{\tilde{\mu}(CA)}{\mu(C)}.$$

If we let $\epsilon \downarrow 0$, then the left-hand and right-hand side converge to $\tilde{\mu}(B)/\mu(B)$ and $\tilde{\mu}(C)/\mu(C)$, respectively, because for arbitrary $U, V \in \mathcal{U}$ with $B \subset U$ and $C \subset V$, $BA \subset U$ and $CA \subset V$ for sufficiently small $\epsilon > 0$. This shows that $\tilde{\mu}/\mu$ is constant on \mathcal{C}_* . \square

By means of Exercise 5.1 (a), one can easily deduce from Theorem 5.6 that there exists a *right Haar measure* μ on \mathcal{X} , that means, Condition 5.1 in Theorem 5.6 can be replaced with

$$(5.9) \quad \mu(Bx) = \mu(B) \quad \text{for all } x \in \mathcal{X} \text{ and } B \in \mathcal{B}.$$

Note that in case of a commutative group, there is no difference between left and right Haar measures, so we talk about *Haar measures*.

An obvious question is how left and right Haar measures are related. The next result is a first step to clarify this.

Theorem 5.8. *There exists a unique function $J : \mathcal{X} \rightarrow (0, \infty)$ such that for any left Haar measure μ on \mathcal{X} ,*

$$\mu(Bx) = J(x)\mu(B) \quad \text{for all } x \in \mathcal{X} \text{ and } B \in \mathcal{B}.$$

Moreover, J is continuous, and

$$J(xy) = J(x)J(y) \quad \text{for arbitrary } x, y \in \mathcal{X}.$$

Here is an immediate consequence of this result:

Corollary 5.9. *Suppose that \mathcal{X} is a compact topological group. Then any left Haar measure on \mathcal{X} is also a right Haar measure.*

The reason is that $J(x^z) = J(x)^z$ for arbitrary $x \in \mathcal{X}$ and integers z . Consequently, if $J(x) \neq 1$ for some $x \in \mathcal{X}$, then J is unbounded on \mathcal{X} . But a continuous function on a compact set is bounded, whence $J \equiv 1$ in case of a compact topological group \mathcal{X} .

Proof of Theorem 5.8. If μ is a left Haar measure on \mathcal{X} , then for fixed $x \in \mathcal{X}$,

$$B \mapsto \mu(Bx)$$

is easily seen to define a left Haar measure on \mathcal{X} , too. But this means that there exists a unique constant $J(x) > 0$ such that

$$\mu(Bx) = J(x)\mu(B) \quad \text{for all } B \in \mathcal{B}.$$

For $x, y \in \mathcal{X}$ and $B \in \mathcal{B}$,

$$J(xy)\mu(B) = \mu(Bxy) = J(y)\mu(Bx) = J(y)J(x)\mu(B).$$

Since $0 < \mu(B) < \infty$ in case of B being compact with non-empty interior, this proves the equation $J(xy) = J(x)J(y)$ for $x, y \in \mathcal{X}$.

To show continuity of J , let $A_\epsilon := B_\epsilon \cap B_\epsilon^{-1}$ with $B_\epsilon = \{x \in \mathcal{X} : d(x, e) \leq \epsilon\}$. Further, let $C \in \mathcal{C}$ with non-empty interior. Then for any $\kappa > 1$ there exists an $\epsilon > 0$ such that $\mu(CA_\epsilon) \leq \mu(C)(1 + \kappa)$. Thus for $x \in A_\epsilon$,

$$J(x) = \frac{\mu(Cx)}{\mu(C)} \leq \frac{\mu(CA_\epsilon)}{\mu(C)} \leq 1 + \kappa$$

and

$$J(x) = \frac{1}{J(x^{-1})} = \frac{\mu(C)}{\mu(Cx^{-1})} \geq \frac{\mu(C)}{\mu(CA_\epsilon)} \geq \frac{1}{1 + \kappa}.$$

This proves continuity of J at e , and continuity at arbitrary points in \mathcal{X} is easily deduced from the equation $J(xy) = J(x)J(y)$ for arbitrary $x, y \in \mathcal{X}$. \square

The function J in Theorem 5.8 is called the *modular function* of the locally compact topological group \mathcal{X} . It yields a Radon–Nikodym derivative of a right-invariant with respect to a left-invariant measure.

Theorem 5.10. *Let μ be a left Haar measure on \mathcal{X} . Then*

$$\tilde{\mu}(B) := \int_B J(x)^{-1} \mu(dx)$$

defines a right-invariant measure on \mathcal{X} . That means, $\tilde{\mu}$ satisfies the same properties as μ , except that (5.1) is replaced with (5.9).

As a preparation for the proof of Theorem 5.10 the reader should verify the following two facts about integrals with respect to left-invariant measures.

Exercise 5.11. Let μ be a left Haar measure on \mathcal{X} . Show that for arbitrary $y \in \mathcal{X}$ and measurable functions $h : \mathcal{X} \rightarrow \overline{\mathbb{R}}$,

$$\int h(yx) \mu(dx) = \int h d\mu \quad \text{and} \quad \int h(xy) \mu(dx) = J(y)^{-1} \int h d\mu,$$

provided that $\int h d\mu$ is well-defined in $\overline{\mathbb{R}}$.

Proof of Theorem 5.10. We start with an arbitrary measure μ satisfying the regularity conditions (5.2) and (5.3) and consider any measure $\tilde{\mu}$ with a continuous density $f = d\tilde{\mu}/d\mu > 0$.

Since $f > 0$, $\tilde{\mu}(B) > 0$ whenever $\mu(B) > 0$. Hence, $\tilde{\mu}$ satisfies property (5.5) whenever μ does. Note also that for $C \in \mathcal{C}$,

$$\tilde{\mu}(C) \leq \max_{x \in C} f(x) \mu(C),$$

so $\tilde{\mu}$ satisfies property (5.4) whenever μ does.

To show that $\tilde{\mu}$ satisfies (5.3), note first that for any $U \in \mathcal{U}$,

$$\tilde{\mu}(U) = \lim_{R \rightarrow \infty} \tilde{\mu}(U \cap \{R^{-1} < f < R\}).$$

Since each set $\{R^{-1} < f < R\}$ is open itself, it suffices to consider open sets U such that $a \leq f \leq b$ on U with $0 < a < b$. In particular, $a\mu(B) \leq \tilde{\mu}(B) \leq b\mu(B)$ for all measurable sets $B \subset U$. By assumption, there exist compact subsets $C_1 \subset C_2 \subset C_3 \subset \dots$ of U such that $\mu(C_n) \rightarrow \mu(U)$ as $n \rightarrow \infty$. Now one can easily show that $\tilde{\mu}(C_n) \rightarrow \tilde{\mu}(U)$ as $n \rightarrow \infty$, distinguishing the cases $\mu(U) = \infty$ and $\mu(U) < \infty$. Hence $\tilde{\mu}$ satisfies (5.3), too.

To show that $\tilde{\mu}$ satisfies (5.2), we fix an arbitrary constant $\gamma > 1$. Then any set $B \in \mathcal{B}$ may be decomposed as

$$B = \bigcup_{z \in \mathbb{Z}} B_z \quad \text{with} \quad B_z := B \cap \{\gamma^z \leq f < \gamma^{z+1}\}.$$

By assumption, for any integer z there exists an open set U_z such that $B_z \subset U_z \subset \{f < \gamma^{z+1}\}$ and $\mu(U_z) \leq \gamma\mu(B_z)$. Consequently, $U := \bigcup_{z \in \mathbb{Z}} U_z$ is an open set containing B such that

$$\tilde{\mu}(U) \leq \sum_{z \in \mathbb{Z}} \tilde{\mu}(U_z) \leq \sum_{z \in \mathbb{Z}} \gamma^{z+1} \mu(U_z) \leq \sum_{z \in \mathbb{Z}} \gamma^{z+2} \mu(B_z) \leq \sum_{z \in \mathbb{Z}} \gamma^2 \tilde{\mu}(B_z) = \gamma^2 \tilde{\mu}(B).$$

Hence, $\tilde{\mu}$ satisfies (5.2), too.

It remains to verify (5.9) for $\tilde{\mu}$ in case of $f(x) = J(x)^{-1}$ and μ being a left Haar measure on \mathcal{X} . Indeed, for measurable functions $h : \mathcal{X} \rightarrow [0, \infty)$ and $y \in \mathcal{X}$,

$$\begin{aligned} \int h(xy) \tilde{\mu}(dx) &= \int h(xy) J(x)^{-1} \mu(dx) \\ &= J(y) \int h(xy) J(xy)^{-1} \mu(dx) \\ &= \int h(x) J(x)^{-1} \mu(dx) \\ &= \int h(x) \tilde{\mu}(dx), \end{aligned}$$

where the second last step follows from Exercise 5.11. □

5.3 Some Explicit Constructions

In this section we present a few examples of Haar measure.

Discrete groups. Let (\mathcal{X}, \cdot) be an arbitrary group. With the metric $d(x, y) := 1_{[x \neq y]}$, the topology \mathcal{U} and the Borel σ -field \mathcal{B} coincide with $\mathcal{P}(\mathcal{X})$, and \mathcal{C} is the family of finite subsets of \mathcal{X} . Thus (\mathcal{X}, \cdot, d) is a locally compact topological group. Here, the counting measure μ ,

$$\mu(B) := \#B,$$

is a left and right Haar measure on \mathcal{X} . This follows essentially from the fact that for any fixed $x_o \in \mathcal{X}$, the mappings $x \mapsto x_o x$ and $x \mapsto x x_o$ are bijective from \mathcal{X} to \mathcal{X} , see Exercise 5.1 (b).

Lebesgue measure on \mathbb{R}^d . Lebesgue measure Leb on (the Borel subsets of) \mathbb{R}^d satisfies

$$\text{Leb}(x + B) = \text{Leb}(B)$$

for arbitrary $x \in \mathbb{R}^d$ and $B \in \text{Borel}(\mathbb{R}^d)$. Consequently, it is the unique Haar measure μ on the additive group \mathbb{R}^d such that $\mu([0, 1]^d) = 1$.

Haar measure on linear groups. Let (\mathcal{X}, \cdot) be the set of nonsingular matrices $\mathbf{x} \in \mathbb{R}^{d \times d}$ with matrix multiplication. As mentioned before, \mathcal{X} is an open subset of $\mathbb{R}^{d \times d}$. For any fixed $\mathbf{x}_o \in \mathcal{X}$,

$$L_{\mathbf{x}_o}(\mathbf{x}) := \mathbf{x}_o \mathbf{x} \quad \text{and} \quad R_{\mathbf{x}_o}(\mathbf{x}) := \mathbf{x} \mathbf{x}_o$$

define bijective mappings from \mathcal{X} to \mathcal{X} . They may also be viewed as bijective linear mappings from $\mathbb{R}^{d \times d}$ to $\mathbb{R}^{d \times d}$. If we identify $\mathbb{R}^{d \times d}$ with \mathbb{R}^{d^2} in the usual fashion, then the determinant of these linear mappings equals

$$\det(L_{\mathbf{x}_o}) = \det(R_{\mathbf{x}_o}) = \det(\mathbf{x}_o)^d,$$

see Exercise 5.12. With Lebesgue measure Leb on $\mathbb{R}^{d \times d}$, this means that for any measurable function $h : \mathcal{X} \rightarrow [0, \infty)$,

$$\int_{\mathcal{X}} h(\mathbf{x}_o \mathbf{x}) d\mathbf{x} = \int_{\mathcal{X}} h(\mathbf{x} \mathbf{x}_o) d\mathbf{x} = |\det(\mathbf{x}_o)|^{-d} \int_{\mathcal{X}} h(\mathbf{x}) d\mathbf{x},$$

where $d\mathbf{x}$ stands for $\text{Leb}(d\mathbf{x})$. This implies that

$$\mu(B) := \int_B |\det(\mathbf{x})|^{-d} d\mathbf{x}$$

defines a left and right Haar measure on \mathcal{X} . For if $\mathbf{x}_o \in \mathcal{X}$, then

$$\begin{aligned} \mu(\mathbf{x}_o^{-1} B) &= \int_{\mathcal{X}} 1_B(\mathbf{x}_o \mathbf{x}) |\det(\mathbf{x})|^{-d} d\mathbf{x} \\ &= |\det(\mathbf{x}_o)|^d \int_{\mathcal{X}} 1_B(\mathbf{x}_o \mathbf{x}) |\det(\mathbf{x}_o \mathbf{x})|^{-d} d\mathbf{x} \\ &= \int_{\mathcal{X}} 1_B(\mathbf{x}) |\det(\mathbf{x})|^{-d} d\mathbf{x} \\ &= \mu(B), \end{aligned}$$

and analogously one can show that $\mu(B \mathbf{x}_o^{-1}) = \mu(B)$.

Exercise 5.12. Let $\text{vec} : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d^2}$ be given by

$$\text{vec}(\mathbf{x}) := (\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_d^\top)^\top$$

for $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$ with columns $\mathbf{x}_j \in \mathbb{R}^d$. Show that for any fixed $\mathbf{x}_o \in \mathcal{X}$ and arbitrary $\mathbf{x} \in \mathbb{R}^{d \times d}$,

$$\text{vec}(\mathbf{x}_o \mathbf{x}) = \mathbf{L}_{\mathbf{x}_o} \text{vec}(\mathbf{x}) \quad \text{and} \quad \text{vec}(\mathbf{x} \mathbf{x}_o) = \mathbf{R}_{\mathbf{x}_o} \text{vec}(\mathbf{x})$$

with matrices $\mathbf{L}_{\mathbf{x}_o}, \mathbf{R}_{\mathbf{x}_o} \in \mathbb{R}^{d^2 \times d^2}$ such that

$$\det(\mathbf{L}_{\mathbf{x}_o}) = \det(\mathbf{x}_o)^d = \det(\mathbf{R}_{\mathbf{x}_o}).$$

A general consideration about compact topological groups. Suppose that \mathcal{X} is a compact topological group, and suppose that P is a left-invariant probability measure on \mathcal{X} , that means,

$$P(xB) = P(B) \quad \text{for all } x \in \mathcal{X} \text{ and } B \in \mathcal{B}.$$

Then P is also right-invariant and inversion-invariant, that means,

$$P(xB) = P(Bx) = P(B^{-1}) = P(B) \quad \text{for all } x \in \mathcal{X} \text{ and } B \in \mathcal{B}.$$

If \tilde{P} is another left-invariant probability measure on \mathcal{X} , then $\tilde{P} \equiv P$.

Proof. Let X and Y be stochastically independent random variables on some probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in $(\mathcal{X}, \mathcal{B})$, where $\mathbb{P}^X = P$ while \mathbb{P}^Y is an arbitrary distribution Q . Then $Y^{-1}X$ has distribution P , too, because for any $B \in \mathcal{B}$,

$$\mathbb{P}(Y^{-1}X \in B) = \mathbb{P}(X \in YB) = \underbrace{\mathbb{E} \mathbb{P}(X \in YB \mid \sigma(Y))}_{=P(YB)=P(B) \text{ a.s.}} = P(B).$$

In particular, if $Q = P$, then for any $B \in \mathcal{B}$,

$$P(B) = \mathbb{P}(Y^{-1}X \in B) = \mathbb{P}(X^{-1}Y \in B^{-1}) = \mathbb{P}(Y \in XB^{-1}) = P(B^{-1}).$$

But this implies right-invariance of P , because for all $x \in \mathcal{X}$ and $B \in \mathcal{B}$,

$$P(Bx) = P((Bx)^{-1}) = P(x^{-1}B^{-1}) = P(B^{-1}) = P(B).$$

Finally, if Q is left-invariant itself, then for all $B \in \mathcal{B}$,

$$Q(B) = \mathbb{P}(X^{-1}Y \in B) = \mathbb{P}(Y^{-1}X \in B^{-1}) = P(B^{-1}) = P(B).$$

□

Haar measure on orthogonal groups. Let \mathcal{X} and $\mathcal{X}_{\text{orth}}$ be the groups of invertible and orthogonal matrices in $\mathbb{R}^{d \times d}$, respectively, the binary operation being matrix multiplication. In view of the previous consideration, it suffices to construct a random variable \mathbf{X} with values in $\mathcal{X}_{\text{orth}}$ such that

$$\mathbb{P}(\mathbf{X} \in \mathbf{x}B) = \mathbb{P}(\mathbf{X} \in B) \quad \text{for all } \mathbf{x} \in \mathcal{X}_{\text{orth}} \text{ and } B \in \mathcal{B}.$$

Then the distribution $P := \mathbb{P}^{\mathbf{X}}$ is the unique left Haar probability measure on $\mathcal{X}_{\text{orth}}$, and it happens to be left-, right- and inversion-invariant.

A particular construction starts with a random matrix $\mathbf{Z} = (Z_{ij})_{i,j=1}^d$ with d^2 stochastically independent random variables $Z_{ij} \sim \mathcal{N}(0, 1)$. Writing $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_d]$, the columns \mathbf{Z}_j are stochastically independent with standard Gaussian distribution $\mathcal{N}_d(\mathbf{0}, \mathbf{e})$. For $2 \leq k \leq d$,

$$\mathbb{P}(\mathbf{Z}_k \in \text{span}(\mathbf{Z}_1, \dots, \mathbf{Z}_{k-1})) = \mathbb{E} \mathbb{P}(\mathbf{Z}_k \in \text{span}(\mathbf{Z}_1, \dots, \mathbf{Z}_{k-1}) \mid \mathbf{Z}_1, \dots, \mathbf{Z}_{k-1}) = 0,$$

because $\mathbb{P}(\mathbf{Z}_k \in \mathbb{V}) = 0$ for any fixed linear space $\mathbb{V} \subset \mathbb{R}^d$ with $\dim(\mathbb{V}) < d$. Consequently, the matrix \mathbf{Z} is invertible almost surely. This allows us to apply the mapping $\Psi : \mathcal{X} \rightarrow \mathcal{X}_{\text{orth}}$, $\Psi(\mathbf{z}) := \mathbf{z}(\mathbf{z}^\top \mathbf{z})^{-1/2}$. Defining $\Psi(\mathbf{z}) := \mathbf{e}$ for singular matrices $\mathbf{z} \in \mathbb{R}^{d \times d}$, the distribution P of the random matrix

$$\mathbf{X} := \Psi(\mathbf{Z})$$

is left-invariant. Indeed, it is well-known from probability theory that the distribution of \mathbf{Z} does not change if we multiply it from the left with an arbitrary matrix $\mathbf{u} \in \mathcal{X}_{\text{orth}}$. Hence, P is also the distribution of

$$\Psi(\mathbf{u}\mathbf{Z}) = \mathbf{u}\mathbf{Z}(\mathbf{Z}^\top \mathbf{u}^\top \mathbf{u}\mathbf{Z})^{-1/2} = \mathbf{u}\mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1/2} = \mathbf{u}\mathbf{X}.$$

Another explicit construction would be to apply the Gram–Schmidt orthogonalization procedure to the columns of \mathbf{Z} . That means, we define $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_d]$ via

$$\mathbf{X}_1 := \|\mathbf{Z}_1\|^{-1} \mathbf{Z}_1$$

and inductively

$$\mathbf{X}_k := \left\| \mathbf{Z}_k - \sum_{i=1}^{k-1} \mathbf{X}_i \mathbf{X}_i^\top \mathbf{Z}_k \right\|^{-1} \left(\mathbf{Z}_k - \sum_{i=1}^{k-1} \mathbf{X}_i \mathbf{X}_i^\top \mathbf{Z}_k \right)$$

for $k = 2, \dots, d$. Again, one can show that replacing \mathbf{Z} with $\mathbf{u}\mathbf{Z}$ for some $\mathbf{u} \in \mathcal{X}_{\text{orth}}$ results in replacing \mathbf{X} with $\mathbf{u}\mathbf{X}$.

Note that the construction via Gram–Schmidt implies that for a random matrix $\mathbf{X} \in \mathcal{X}_{\text{orth}}$ with left-invariant distribution, any column is uniformly distributed on the unit sphere of \mathbb{R}^d (with respect to the standard Euclidean norm).

Part II

Mathematical Statistics

Chapter 6

Measurement Series and Estimators of Location

6.1 Statistical Experiments and Point Estimators

Before we discuss estimation of a location parameter, let us introduce some general terminology.

Definition 6.1 (Statistical experiment). A *statistical experiment* is a triplet $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ consisting of a measurable space (Ω, \mathcal{A}) , the *sample space*, and a family of probability distributions \mathbb{P}_θ on (Ω, \mathcal{A}) , depending on a *parameter* θ in a *parameter space* Θ .

The sample space (Ω, \mathcal{A}) represents all possible data sets one could observe. The elements of the parameter space Θ represent potential values of an *unknown true parameter* $\theta \in \Theta$. We assume that the observed data are a (realization of a) random variable with distribution \mathbb{P}_θ . In what follows, the symbol θ may denote this particular true parameter or a potential parameter. It should become clear from the context in which sense θ is meant. The dependency of probabilities, expectations, variances etc. on θ will be denoted by a corresponding subscript, leading to \mathbb{P}_θ , \mathbb{E}_θ , Var_θ etc.

Of course there is some redundancy in the definition of a statistical experiment \mathcal{E} , because specifying the family $(\mathbb{P}_\theta)_{\theta \in \Theta}$ implies the specification of the measurable space (Ω, \mathcal{A}) . But sometimes we shall replace \mathcal{A} with certain sub- σ -fields, so the current definition is useful.

Sometimes one is interested in a function $g(\theta)$ of the true parameter θ with values in some metric space (\mathbb{G}, d) , where $g : \Theta \rightarrow \mathbb{G}$ is given. In the simplest case, one would like to deduce from observed data $\omega \in \Omega$ a simple guess $\hat{g}(\omega) \in \mathbb{G}$ of $g(\theta)$.

Definition 6.2 (Point estimator). A *point estimator of $g(\theta)$* (short: an *estimator*) is a measurable¹ function

$$\hat{g} : \Omega \rightarrow \mathbb{G}.$$

To compare different estimators \hat{g} , one can quantify their inaccuracy for instance by their mean

¹ \mathcal{A} -Borel(\mathbb{R}^q)-measurable

squared error,

$$\mathbb{E}_\theta(d(\hat{g}, g(\theta))^2).$$

Note that this quantity depends on the parameter θ . In general, it may happen that one estimator \hat{g}_1 is strictly better than another estimator \hat{g}_2 in a certain region of the parameter space Θ but strictly worse in another region.

6.2 Estimators of Location

Simple location families. Let $\theta \in \mathbb{R}$ be an unknown parameter which has to be estimated in an experiment which yields n single measurements X_1, \dots, X_n . Suppose that

$$X_i = \theta + \epsilon_i, \quad 1 \leq i \leq n,$$

with independent random measurement errors $\epsilon_1, \dots, \epsilon_n$ with a known distribution P_0 . This leads to the statistical experiment

$$\mathcal{E} = (\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (\mathbb{P}_\theta)_{\theta \in \mathbb{R}})$$

with

$$\mathbb{P}_\theta := P_\theta^{\otimes n} \quad \text{and} \quad P_\theta = P_0 \star \delta_\theta.$$

Here ‘ \star ’ denotes convolution², and δ_θ is the Dirac measure at the point θ . In other words, \mathbb{P}_θ describes the distribution of a random vector with independent components having distribution function F_θ , where

$$F_\theta(x) := F_0(x - \theta),$$

and F_0 is the distribution function of the error distribution P_0 . Such a statistical experiment is called a *simple location family*.

In what follows, we write \mathbf{x} instead of ω for a sample in \mathbb{R}^n . The random variables X_i are just the coordinate functions $\mathbf{x} = (x_i)_{i=1}^n \mapsto X_i(\mathbf{x}) := x_i$, and $\mathbf{X} := (X_i)_{i=1}^n$ is the identity function.

Equivariant estimators. For a vector $\mathbf{x} = (x_i)_{i=1}^n \in \mathbb{R}^n$ and any number $a \in \mathbb{R}$ let

$$a + \mathbf{x} := \mathbf{x} + a := (x_i + a)_{i=1}^n$$

The simple location family \mathcal{E} has the property that for arbitrary $\theta, a \in \mathbb{R}$,

$$\mathbf{X} \sim \mathbb{P}_\theta \quad \text{if and only if} \quad \mathbf{X} + a \sim \mathbb{P}_{\theta+a}.$$

This motivates the following property of an estimator $\hat{\theta}$ of θ :

Definition 6.3 (Equivariance). An estimator $\hat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$ is called *equivariant*, if

$$\hat{\theta}(\mathbf{x} + a) = \hat{\theta}(\mathbf{x}) + a \quad \text{for all } \mathbf{x} \in \mathbb{R}^n \text{ and } a \in \mathbb{R}.$$

²For distributions P and Q on the real line, $P \star Q$ denotes the distribution of $X + Y$ with independent random variables $X \sim P$ and $Y \sim Q$. In particular, for $\theta \in \mathbb{R}$, $P \star \delta_\theta$ is the distribution of $X + \theta$ with $X \sim P$.

Note that the sample mean,

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i,$$

as well as the sample median are equivariant estimators. Concerning the sample mean, it follows from the weak law of large numbers that for a fixed distribution P_0 with mean $\int x P_0(dx) = 0$,

$$\lim_{n \rightarrow \infty} \mathbb{E}_\theta |\bar{X} - \theta| = 0.$$

If in addition, P_0 has finite variance σ^2 , then even

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

Risk functions. For an arbitrary estimator $\hat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$, we consider its risk function $R(\hat{\theta}, \cdot) : \mathbb{R} \rightarrow [0, \infty]$ with

$$R(\hat{\theta}, \theta) := \mathbb{E}_\theta ((\hat{\theta} - \theta)^2) = \int_{\mathbb{R}^n} (\hat{\theta} - \theta)^2 d\mathbb{P}_\theta,$$

the mean squared error of $\hat{\theta}$ in case of the true parameter being θ .

In case of an equivariant estimator $\hat{\theta}$, its risk function is constant: For arbitrary $\theta \in \mathbb{R}$,

$$R(\hat{\theta}, \theta) = R(\hat{\theta}) := R(\hat{\theta}, 0) = \mathbb{E}_0(\hat{\theta}^2).$$

More generally, if $\hat{\theta}$ is equivariant, then for any measurable function $h : \mathbb{R} \rightarrow \mathbb{R}$ and arbitrary $\theta \in \mathbb{R}$,

$$\mathbb{E}_\theta h(\hat{\theta} - \theta) = \mathbb{E}_0 h(\hat{\theta}),$$

provided that $\mathbb{E}_0 h(\hat{\theta})$ is well-defined.

6.3 Constructing an Optimal Equivariant Estimator

An equivariant estimator $\hat{\theta}_*$ is called *optimal (among all equivariant estimators)* if

$$R(\hat{\theta}_*) \leq R(\hat{\theta}) \quad \text{for any equivariant estimator } \hat{\theta}.$$

If $\mathbf{X} \sim \mathbb{P}_\theta$, then $\mathbf{X} = \theta + \epsilon$ with $\epsilon \sim \mathbb{P}_0$, and for any equivariant estimator $\hat{\theta}$,

$$\hat{\theta}(\mathbf{X}) = \theta + \hat{\theta}(\epsilon).$$

Of course, we don't know ϵ , but at least we know $\mathbf{T}(\epsilon)$ for the particular function $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by

$$\mathbf{T}(\mathbf{x}) := \mathbf{x} - x_1 = (0, x_2 - x_1, \dots, x_n - x_1)^\top.$$

Indeed, this function \mathbf{T} is *invariant* in the sense that

$$\mathbf{T}(\mathbf{x} + a) = \mathbf{T}(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathbb{R}^n \text{ and } a \in \mathbb{R}.$$

Hence, if we observe $\mathbf{X} = \theta + \epsilon$, then

$$\mathbf{T}(\mathbf{X}) = \mathbf{T}(\epsilon).$$

That means, we know at least $\mathbf{T}(\epsilon)$. So we could try to improve the estimator $\widehat{\theta}(\mathbf{X}) = \theta + \widehat{\theta}(\epsilon)$ by subtracting the conditional expectation of $\widehat{\theta}(\epsilon)$, given that $\mathbf{T}(\epsilon)$ is equal to the observed $\mathbf{T}(\mathbf{X})$. That means, we subtract a reasonable guess of $\widehat{\theta}(\epsilon)$ from $\widehat{\theta}(\mathbf{X})$. This idea leads to an optimal equivariant estimator indeed.

Theorem 6.4 (Pitman's improvement). *Let $\widehat{\theta}$ be an equivariant estimator with finite risk $R(\widehat{\theta})$. Then*

$$\widehat{\theta}_* := \widehat{\theta} - \mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T}))$$

defines an optimal equivariant estimator. It is unique in the sense that for any equivariant estimator $\tilde{\theta}$ and arbitrary $\theta \in \mathbb{R}$,

$$R(\tilde{\theta}) = R(\widehat{\theta}_*) \quad \text{implies that} \quad \tilde{\theta} = \widehat{\theta}_* \quad \mathbb{P}_\theta\text{-almost surely.}$$

Remark 6.5 (Invariance and the choice of $\mathbf{T}(\cdot)$). Our particular choice of $\mathbf{T}(\cdot)$ is somewhat arbitrary. In principle one could take any equivariant estimator $\tilde{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$ and define $\mathbf{T}(\mathbf{x}) := \mathbf{x} - \tilde{\theta}(\mathbf{x})$. Inspecting the proof of Theorem 6.4 carefully reveals that Theorem 6.4 remains valid with this definition \mathbf{T} . Our particular version corresponds to $\tilde{\theta} = X_1$ and is convenient for explicit calculations.

Exercise 6.6. Consider an estimator $\tilde{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$. Show that $\tilde{\theta}$ is equivariant if and only if $\mathbf{T}(\mathbf{x}) := \mathbf{x} - \tilde{\theta}(\mathbf{x})$ is invariant, i.e. $\mathbf{T}(\mathbf{x} + a) = \mathbf{T}(\mathbf{x})$ for arbitrary $\mathbf{x} \in \mathbb{R}^n$ and $a \in \mathbb{R}$.

Remark 6.7 (Characterization of optimality). The particular construction in Theorem 6.4 implies that an equivariant estimator $\widehat{\theta}$ with finite risk $R(\widehat{\theta})$ is optimal if and only if

$$\mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T})) = 0 \quad \mathbb{P}_0\text{-almost surely.}$$

Exercise 6.8. Let $\widehat{\theta}_* : \mathbb{R}^n \rightarrow \mathbb{R}$ be an optimal equivariant estimator for θ . Show that $\widehat{\theta}_*$ is unbiased, that means,

$$\mathbb{E}_\theta(\widehat{\theta}_*) = \theta \quad \text{for all } \theta \in \mathbb{R}.$$

Exercise 6.9. Suppose that P_0 is the Laplace distribution on $\{0, 1\}$.

- (a) Before starting to apply the general theory, how would you estimate θ ?
- (b) Determine the conditional distribution of \mathbf{X} , given that $\mathbf{T} = \mathbf{y}$, in case of $\theta = 0$. (Which vectors $\mathbf{y} \in \mathbb{R}^n$ are relevant?)
- (c) Determine the optimal (in terms of mean squared error) equivariant estimator of θ .

Proof of Theorem 6.4. The general theory in Sections 3.2 and 3.3 shows that $\mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T})) = g_*(\mathbf{T})$ with a measurable function $g_* : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$\mathbb{E}_0((\widehat{\theta} - g(\mathbf{T}))^2) = \mathbb{E}_0((\widehat{\theta} - g_*(\mathbf{T}))^2) + \mathbb{E}_0((g(\mathbf{T}) - g_*(\mathbf{T}))^2)$$

for any measurable function $g : \mathbb{R}^n \rightarrow \mathbb{R}$. In particular, for $g \equiv 0$ we obtain the formula

$$(6.1) \quad \mathbb{E}_0(\widehat{\theta}^2) = \mathbb{E}_0((\widehat{\theta} - g_*(\mathbf{T}))^2) + \mathbb{E}_0(g_*(\mathbf{T})^2).$$

Since $\widehat{\theta}$ is equivariant and \mathbf{T} is invariant, $\widehat{\theta}_* = \widehat{\theta} - g_*(\mathbf{T})$ is an equivariant estimator, too: For arbitrary $\mathbf{x} \in \mathbb{R}^n$ and $a \in \mathbb{R}$,

$$\widehat{\theta}_*(\mathbf{x} + a) = \underbrace{\widehat{\theta}(\mathbf{x} + a)}_{=\widehat{\theta}(\mathbf{x})+a} - \underbrace{g_*(\mathbf{T}(\mathbf{x} + a))}_{=g_*(\mathbf{T}(\mathbf{x}))} = \widehat{\theta}(\mathbf{x}) + a - g_*(\mathbf{T}(\mathbf{x})) = \widehat{\theta}_*(\mathbf{x}) + a.$$

Hence, we may rewrite (6.1) as

$$R(\widehat{\theta}) = R(\widehat{\theta}_*) + \mathbb{E}_0(g_*(\mathbf{T})^2).$$

Consequently, $R(\widehat{\theta}) \geq R(\widehat{\theta}_*)$ with equality if and only if $\mathbb{E}_0(g_*(\mathbf{T})^2) = 0$, and this is equivalent to

$$\mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T})) = 0 \quad \mathbb{P}_0\text{-almost surely.}$$

Finally, let $\tilde{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$ be another equivariant estimator with finite risk $R(\tilde{\theta})$. Then $h := \tilde{\theta} - \widehat{\theta}$ is invariant in the sense that $h(\mathbf{x}) = h(\mathbf{T}(\mathbf{x}))$ for arbitrary $\mathbf{x} \in \mathbb{R}^n$. Consequently, if we apply Pitman's recipe to $\tilde{\theta}$ instead of $\widehat{\theta}$, we obtain the estimator

$$\begin{aligned} \tilde{\theta}_* &:= \tilde{\theta} - \mathbb{E}_0(\tilde{\theta} | \sigma(\mathbf{T})) \\ &= \widehat{\theta} + h(\mathbf{T}) - \mathbb{E}_0(\widehat{\theta} + h(\mathbf{T}) | \sigma(\mathbf{T})) \\ &= \widehat{\theta} + h(\mathbf{T}) - \mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T})) - \underbrace{\mathbb{E}_0(h(\mathbf{T}) | \sigma(\mathbf{T}))}_{=h(\mathbf{T}) \text{ a.s.}} \\ &= \widehat{\theta} - \mathbb{E}_0(\widehat{\theta} | \sigma(\mathbf{T})) = \widehat{\theta}_* \quad \mathbb{P}_0\text{-almost surely.} \end{aligned}$$

Moreover, since $\tilde{\theta}_* - \widehat{\theta}_*$ is invariant, $\mathbb{P}_\theta(\tilde{\theta}_* \neq \widehat{\theta}_*) = \mathbb{P}_0(\tilde{\theta}_* \neq \widehat{\theta}_*) = 0$ for arbitrary $\theta \in \mathbb{R}$. \square

In case of P_0 having a density with respect to Lebesgue measure, there is an explicit formula for the optimal equivariant estimator $\widehat{\theta}_*$:

Corollary 6.10. *Suppose that P_0 has a density f_0 with respect to Lebesgue measure on \mathbb{R} , and suppose that there exists an equivariant estimator with finite risk. Then there exists a Borel set $B_* \subset \mathbb{R}^n$ such that $\mathbb{P}_0(\mathbf{T}^{-1}(B_*)) = 1$,³ and for each $\mathbf{x} \in \mathbf{T}^{-1}(B_*)$, the optimal equivariant estimator $\widehat{\theta}_*$ is given by*

$$\widehat{\theta}_*(\mathbf{x}) = \int_{\mathbb{R}} \theta f_\theta(\mathbf{x}) d\theta / \int_{\mathbb{R}} f_\theta(\mathbf{x}) d\theta,$$

where $f_\theta(\mathbf{x}) := \prod_{i=1}^n f_\theta(x_i)$ and $f_\theta(x) := f_0(x - \theta)$ for real numbers x . In other words, $\widehat{\theta}_*(\mathbf{x})$ is the mean of the probability distribution $Q_{\mathbf{x}}$ on \mathbb{R} with density

$$\theta \mapsto f_\theta(\mathbf{x}) / \int_{\mathbb{R}} f_\eta(\mathbf{x}) d\eta.$$

³ $\mathbb{P}_\theta(\mathbf{T}^{-1}(B_*)) = 1$ for all $\theta \in \mathbb{R}$ by invariance of \mathbf{T} .

Example 6.11 (Gaussian distributions). Suppose that $P_0 = \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. Then

$$\hat{\theta}_* = \bar{X}.$$

This follows from Corollary 6.10 and the following calculations: $f_\theta(x) = C \exp(-x^2/(2\sigma^2))$ with $C = (2\pi\sigma^2)^{-1/2}$, whence

$$f_\theta(\mathbf{x}) = C^n \exp\left(-\frac{\|\mathbf{x} - \theta\|^2}{2\sigma^2}\right).$$

But

$$\|\mathbf{x} - \theta\|^2 = \|\mathbf{x} - \bar{x}\|^2 + n(\theta - \bar{x})^2$$

with $\bar{x} := n^{-1} \sum_{i=1}^n x_i$, so

$$f_\theta(\mathbf{x}) = f_{\bar{x}}(\mathbf{x}) \exp\left(-\frac{(\theta - \bar{x})^2}{2\sigma^2/n}\right),$$

and this implies that the distribution $Q_{\mathbf{x}}$ in Corollary 6.10 is equal to $\mathcal{N}(\bar{x}, \sigma^2/n)$. Hence

$$\hat{\theta}_*(\mathbf{x}) = \text{mean}(\mathcal{N}(\bar{x}, \sigma^2/n)) = \bar{x}.$$

Example 6.12 (Uniform distributions). Suppose that $P_0 = \text{Unif}([-\sigma, \sigma])$ for some $\sigma > 0$. Then

$$\hat{\theta}_*(\mathbf{x}) = (\min(\mathbf{x}) + \max(\mathbf{x}))/2$$

with $\min(\mathbf{x})$ and $\max(\mathbf{x})$ denoting the minimum and maximum of $\{x_1, \dots, x_n\}$, respectively. This follows from the following considerations: Since $f_0(x) = (2\sigma)^{-1} 1_{[-\sigma \leq x \leq \sigma]}$,

$$\begin{aligned} f_\theta(\mathbf{x}) &= (2\sigma)^{-n} \prod_{i=1}^n 1_{[-\sigma \leq x_i - \theta \leq \sigma]} \\ &= (2\sigma)^{-n} \prod_{i=1}^n 1_{[x_i - \sigma \leq \theta \leq x_i + \sigma]} \\ &= (2\sigma)^{-n} 1_{[\max(\mathbf{x}) - \sigma \leq \theta \leq \min(\mathbf{x}) + \sigma]}. \end{aligned}$$

Hence the distribution $Q_{\mathbf{x}}$ in Corollary 6.10 is the uniform distribution on the interval with endpoints $\max(\mathbf{x}) - \sigma$ and $\min(\mathbf{x}) + \sigma$, unless $\max(\mathbf{x}) - \min(\mathbf{x}) > 2\sigma$. (Note that $\max(\mathbf{X}) - \min(\mathbf{X}) < 2\sigma$ almost surely.) Consequently,

$$\hat{\theta}_*(\mathbf{x}) = \text{midpoint of } [\max(\mathbf{x}) - \sigma, \min(\mathbf{x}) + \sigma] = (\min(\mathbf{x}) + \max(\mathbf{x}))/2.$$

(This definition makes sense no matter how large the difference $\max(\mathbf{x}) - \min(\mathbf{x})$ is.)

Exercise 6.13. Suppose that $P_0 = \text{Unif}[-\sigma, \sigma]$. As shown before, the optimal equivariant estimator of θ is given by $\hat{\theta}_*(\mathbf{x}) = (\min(\mathbf{x}) + \max(\mathbf{x}))/2$.

(a) Determine the risk of \bar{X} .

(b) Show that the risk of $\hat{\theta}_*$ is of order $O(n^{-2})$.

Bonus question: Show that

$$R(\hat{\theta}_*) = \frac{2\sigma^2}{(n+1)(n+2)}.$$

Remark 6.14 (Maximum-likelihood estimation). In case of P_0 having density f_0 , the function

$$\theta \mapsto f_\theta(\mathbf{X}) = \prod_{i=1}^n f_\theta(X_i)$$

is the so-called likelihood function. For any number θ one may interpret $f_\theta(\mathbf{X})$ as a measure of plausibility of θ being equal to the true parameter. Indeed a standard estimator of the true parameter θ would be the maximum-likelihood estimator

$$\hat{\theta}_{\text{ML}} := \arg \max_{\theta \in \mathbb{R}} f_\theta(\mathbf{X}),$$

provided the latter is uniquely defined. Our previous calculations show that $\hat{\theta}_* = \hat{\theta}_{\text{ML}} = \bar{X}$ in case of P_0 being a centered Gaussian distribution.

The higher popularity of $\hat{\theta}_{\text{ML}}$ in comparison with $\hat{\theta}_*$ is due to the fact that the latter estimator is rather difficult to compute explicitly in non-Gaussian models. Moreover, in many settings one can show that

$$f_\theta(\mathbf{X}) \approx f_{\hat{\theta}_{\text{ML}}}(\mathbf{X}) \exp\left(-\frac{(\theta - \hat{\theta}_{\text{ML}})^2}{2\hat{\gamma}^2}\right)$$

for some random variable $\hat{\gamma} > 0$ such that $\hat{\gamma} \rightarrow_p 0$ as $n \rightarrow \infty$. Hence $Q_{\mathbf{X}} \approx \mathcal{N}(\hat{\theta}_{\text{ML}}, \hat{\gamma}^2)$ for large sample sizes n , and $\hat{\theta}_{\text{ML}}$ seems to be a good surrogate for $\hat{\theta}_*$.

Exercise 6.15. Suppose that the error distribution P_0 is the standard exponential distribution. That means, its density is given by

$$f_0(x) = \begin{cases} 0 & \text{if } x < 0, \\ \exp(-x) & \text{if } x \geq 0. \end{cases}$$

- (a) Determine $f_\theta(\mathbf{x})$ for $\theta \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$ in terms of $\min(\mathbf{x})$ and $x_+ := \sum_{i=1}^n x_i$.
- (b) Determine the maximum likelihood estimator $\hat{\theta}_{\text{ML}}$.
- (c) Determine the optimal equivariant estimator $\hat{\theta}_*$.

Proof of Corollary 6.10. Consider the linear transformation given by

$$\mathbf{x} \mapsto (x_1, x_2 - x_1, \dots, x_n - x_1)^\top = \mathbf{A}\mathbf{x}$$

with the lower triangular matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ -1 & 1 & 0 & \dots & 0 \\ -1 & 0 & 1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ -1 & 0 & \dots & 0 & 1 \end{bmatrix}.$$

Its inverse is given by

$$\mathbf{y} \mapsto (y_1, y_2 + y_1, \dots, y_n + y_1)^\top = \mathbf{A}^{-1}\mathbf{y},$$

and $\det(\mathbf{A}) = 1 = \det(\mathbf{A}^{-1})$. Hence, by the transformation formula for Lebesgue integrals,

$$\begin{aligned}\mathbb{P}_0(\mathbf{A}\mathbf{X} \in B) &= \int_{\mathbb{R}^n} 1_{[\mathbf{A}\mathbf{x} \in B]} f_0(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^n} 1_{[\mathbf{y} \in B]} f_0(\mathbf{A}^{-1}\mathbf{y}) d\mathbf{y} \\ &= \int_B f_0(\mathbf{A}^{-1}\mathbf{y}) d\mathbf{y}\end{aligned}$$

for any Borel set $B \subset \mathbb{R}^n$. This shows that the distribution of $\mathbf{A}\mathbf{X}$ under \mathbb{P}_0 has density⁴

$$\mathbf{y} \mapsto f_0(\mathbf{A}^{-1}\mathbf{y}) = f_0(y_1, y_2 + y_1, \dots, y_n + y_1).$$

Note that $\mathbf{A}\mathbf{X} = (X_1, T_2, \dots, T_n)^\top$ while $T_1 \equiv 0$. Thus the considerations in Example 4.10 imply the following: The distribution of (T_2, \dots, T_n) is given by the density

$$(t_2, \dots, t_n) \mapsto g(t_2, \dots, t_n) := \int_{\mathbb{R}} f_0(u, t_2 + u, \dots, t_n + u) du.$$

In particular, there exists a Borel set $B_o \subset \mathbb{R}^n$ with $\mathbb{P}_0(\mathbf{T} \in B_o) = 1$ such that for any $\mathbf{t} \in B_o$,

$$0 < g(t_2, \dots, t_n) := \int_{\mathbb{R}} f_0(u, t_2 + u, \dots, t_n + u) du < \infty,$$

and the conditional distribution of X_1 , given that $\mathbf{T} = \mathbf{t}$, has density

$$u \mapsto g(t_2, \dots, t_n)^{-1} f_0(u, t_2 + u, \dots, t_n + u).$$

Coming back to our estimator $\hat{\theta}$ with finite risk $\mathbb{E}_0(\hat{\theta}^2)$, note that $\hat{\theta}(\mathbf{x}) = x_1 + h(\mathbf{T}(\mathbf{x}))$ for some measurable function h on \mathbb{R}^n , and

$$\infty > \mathbb{E}_0|\hat{\theta}| = \int_{\mathbb{R}^{n-1}} \int_{\mathbb{R}} |u + h(0, t_2, \dots, t_n)| f_0(u, t_2 + u, \dots, t_n + u) du d(t_2, \dots, t_n).$$

Hence for a Borel set $B_* \subset B_o$ with $\mathbb{P}_0(\mathbf{T} \in B_*) = 1$ and arbitrary $\mathbf{t} \in B_*$,

$$\int_{\mathbb{R}} |u| f_0(u, t_2 + u, \dots, t_n + u) du < \infty.$$

In particular, if we plug in $\mathbf{t} = \mathbf{T}(\mathbf{x})$ for some $\mathbf{x} \in \mathbf{T}^{-1}(B_*)$, then

$$f_0(u, t_2 + u, \dots, t_n + u) = f_0(x_1 - x_1 + u, x_2 - x_1 + u, \dots, x_n - x_1 + u) = f_{x_1 - u}(\mathbf{x}),$$

so

$$g(t_2, \dots, t_n) = \int_{\mathbb{R}} f_{x_1 - u}(\mathbf{x}) du = \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta,$$

and

$$\begin{aligned}\mathbb{E}_0(\hat{\theta} | \mathbf{T} = \mathbf{T}(\mathbf{x})) &= \int_{\mathbb{R}} (u + h(\mathbf{T}(\mathbf{x}))) f_{x_1 - u}(\mathbf{x}) du / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta \\ &= x_1 + h(\mathbf{T}(\mathbf{x})) - \int_{\mathbb{R}} (x_1 - u) f_{x_1 - u}(\mathbf{x}) du / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta \\ &= \hat{\theta}(\mathbf{x}) - \int_{\mathbb{R}} \theta f_{\theta}(\mathbf{x}) d\theta / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta.\end{aligned}$$

⁴In this proof, densities are with respect to the corresponding Lebesgue measure.

Consequently,

$$\hat{\theta}_*(\mathbf{x}) = \hat{\theta}(\mathbf{x}) - \mathbb{E}_0(\hat{\theta} | \mathbf{T} = \mathbf{T}(\mathbf{x})) = \int_{\mathbb{R}} \theta f_{\theta}(\mathbf{x}) d\theta / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta$$

for all $\mathbf{x} \in \mathbf{T}^{-1}(B_*)$, as claimed. \square

We end this section with the interesting result that in case of an error distribution P_0 with finite second moment and Lebesgue density f_0 , the optimal equivariant estimator is the sample mean if and only if P_0 is a centered Gaussian distribution.

Theorem 6.16 (Kagan–Linnik–Rao). *Let $n \geq 3$, and let P_0 have finite second moment. Then $\hat{\theta}_* = \bar{X}$ almost surely if and only if $P_0 = \mathcal{N}(0, \sigma^2)$ for some $\sigma \geq 0$.*

Proof of Theorem 6.16. We have verified already that $\hat{\theta}_* = \bar{X}$ in case of P_0 being a centered Gaussian distribution. Hence it suffices to prove the reverse statement.

In what follows all probabilities and expectations refer to the distribution $\mathbb{P} = \mathbb{P}_0$. The proof of Theorem 6.4 shows that optimality of \bar{X} is equivalent to

$$\mathbb{E}(\bar{X} | \sigma(\mathbf{T})) = 0$$

almost surely, where $\mathbf{T} = \mathbf{X} - X_1 = (0, X_2 - X_1, \dots, X_n - X_1)^\top$. In other words,

$$(6.2) \quad \mathbb{E}\left(\sum_{i=1}^n X_i g(\mathbf{T})\right) = 0 \quad \text{whenever } g(\mathbf{T}) \in L^2(\mathbb{P}).$$

With $g \equiv 1$ this implies that $X := X_1$ has mean

$$\mathbb{E}(X) = \int x f_0(x) dx = 0.$$

If $n > 3$, we may take $g(\mathbf{T}) = h(T_2, T_3) = h(X_2 - X_1, X_3 - X_1)$ and deduce from independence of X_1, X_2, \dots, X_n and $\mathbb{E}(X_i) = 0$ that

$$(6.3) \quad \mathbb{E}((X_1 + X_2 + X_3)h(T_2, T_3)) = 0 \quad \text{whenever } h(T_2, T_3) \in L^2(\mathbb{P}).$$

Now let ϕ be the characteristic function of P_0 , i.e.

$$\phi(t) = \mathbb{E}(e^{itX}) = \int e^{itx} P_0(dx)$$

with the imaginary unit $i \in \mathbb{C}$. We know that $\phi : \mathbb{R} \rightarrow \mathbb{C}$ is bounded and continuous with $\phi(0) = 1$. Moreover, since P_0 has finite second moment, ϕ is twice continuously differentiable with derivative $\phi^{(k)}(t) = i^k \mathbb{E}(X^k e^{itX})$ for $k = 1, 2$, so

$$\phi'(0) = i \mathbb{E}(X) = 0 \quad \text{and} \quad \phi''(0) = -\mathbb{E}(X^2).$$

We may apply (6.3) to the (real and imaginary part of the) complex-valued and bounded function $h(z_1, z_2) = e^{is z_1 + it z_2}$ with arbitrary real numbers s, t . This leads to

$$\begin{aligned}
0 &= \mathbb{E}((X_1 + X_2 + X_3)e^{is(X_2 - X_1) + it(X_3 - X_1)}) \\
&= \mathbb{E}(X_1 e^{-i(s+t)X_1} e^{isX_2} e^{itX_3}) \\
&\quad + \mathbb{E}(e^{-i(s+t)X_1} X_2 e^{isX_2} e^{itX_3}) + \mathbb{E}(e^{-i(s+t)X_1} e^{isX_2} X_3 e^{itX_3}) \\
&= \mathbb{E}(X_1 e^{-i(s+t)X_1}) \mathbb{E}(e^{isX_2}) \mathbb{E}(e^{itX_3}) \\
&\quad + \mathbb{E}(e^{-i(s+t)X_1}) \mathbb{E}(X_2 e^{isX_2}) \mathbb{E}(e^{itX_3}) + \mathbb{E}(e^{-i(s+t)X_1}) \mathbb{E}(e^{isX_2}) \mathbb{E}(X_3 e^{itX_3}) \\
&= \phi'(-(s+t))\phi(s)\phi(t) + \phi(-(s+t))\phi'(s)\phi(t) + \phi(-(s+t))\phi(s)\phi'(t),
\end{aligned}$$

where the latter two equalities follow from X_1, X_2, X_3 being independent and identically distributed. Consequently,

$$\phi'(-(s+t))\phi(s)\phi(t) + \phi(-(s+t))(\phi'(s)\phi(t) + \phi(s)\phi'(t)) = 0 \quad \text{for arbitrary } s, t \in \mathbb{R}.$$

Now let

$$c := \max\{t \in (0, \infty] : \phi \neq 0 \text{ on } (-t, t)\}.$$

Then

$$\psi(t) := \frac{\phi'(t)}{\phi(t)}$$

defines a continuous function $\psi : (-c, c) \rightarrow \mathbb{C}$ with $\psi(0) = 0$ and

$$\psi(-(s+t)) + \psi(s) + \psi(t) = 0 \quad \text{whenever } |s|, |t|, |s+t| < c.$$

But this implies that for some $\alpha \in \mathbb{C}$,

$$\psi(t) = \alpha t \quad \text{for all } t \in (-c, c),$$

see Exercise 6.17. In other words,

$$\phi'(t) = \alpha t \phi(t) \quad \text{for } t \in (-c, c).$$

Together with $\phi(0) = 1$, standard results for differential equations imply that

$$\phi(t) = e^{\alpha t^2/2} \quad \text{for } t \in (-c, c).$$

But continuity of ϕ and the definition of c imply that $c = \infty$. For otherwise, continuity of ϕ and the definition of c would imply that $0 = \phi(\pm c) = e^{\alpha c^2/2}$. Consequently,

$$\phi(t) = e^{\alpha t^2/2} \quad \text{for all } t \in \mathbb{R}.$$

Since $\phi''(t) = (\alpha + \alpha^2 t^2)\phi(t)$, we may conclude from $\phi''(0) = \alpha = -\mathbb{E}(X^2)$ that α is a negative real number. It is well-known from probability theory that for any $\sigma \geq 0$, the characteristic function of $\mathcal{N}(\mu, \sigma^2)$ is given by $t \mapsto \exp(it\mu - \sigma^2 t^2/2)$. Hence the characteristic function of P_0 coincides with the characteristic function of $\mathcal{N}(0, \sigma^2)$, where $\sigma := \sqrt{-\alpha}$. Since any probability distribution is uniquely determined by its characteristic distribution, this shows that $P_0 = \mathcal{N}(0, \sigma^2)$. \square

Exercise 6.17. For some $c \in (0, \infty]$, let $\psi : (-c, c) \rightarrow \mathbb{C}$ be a continuous function such that

$$\psi(-(s+t)) + \psi(s) + \psi(t) = 0 \quad \text{whenever } |s|, |t|, |s+t| < c.$$

Show that there exists a constant $\alpha \in \mathbb{C}$ such that

$$\psi(t) = \alpha t \quad \text{for all } t \in (-c, c).$$

Exercise 6.18 (Distribution of order statistics). The contents of this exercise are probably known from other courses in Statistics. Let $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ be the order statistics of independent random variables X_1, \dots, X_n with distribution function F on \mathbb{R} .

(a) Show that for $k \in \{1, 2, \dots, n\}$,

$$\mathbb{P}(X_{(k)} \leq x) = 1 - B_{n, F(x)}(k-1),$$

where $B_{n,p}(\cdot)$ denotes the distribution function of the binomial distribution $\text{Bin}(n, p)$.

(b) Show that for $c \in \{0, 1, \dots, n-1\}$, $B_{n,0}(c) = 1$, $B_{n,1}(c) = 0$ and

$$B_{n,p}(c) = n \binom{n-1}{c} \int_p^1 u^c (1-u)^{n-1-c} du \quad \text{for } p \in [0, 1].$$

(c) Deduce from parts (a-b) that

$$\mathbb{P}(X_{(k)} \leq x) = n \binom{n-1}{k-1} \int_0^{F(x)} u^{k-1} (1-u)^{n-k} du.$$

Exercise 6.19 (Distribution of the sample median). Let X_1, \dots, X_n be independent random variables with density f and differentiable distribution function F on \mathbb{R} . Suppose that $n = 2m + 1$ for some integer $m \geq 1$.

(a) Show that the sample median $M_n := \text{median}(X_1, \dots, X_n)$ has density

$$f_n(x) = n \binom{2m}{m} F(x)^m (1-F(x))^m f(x).$$

(b) Suppose that f is the standard Cauchy density, $f(x) = \pi^{-1}(1+x^2)^{-1}$. For which values of m is $\mathbb{E}(M_n^2) < \infty$?

Remark 1: One can answer (b) without computing $\mathbb{E}(M_n^2)$ explicitly, utilizing rough bounds for $F(x)$.

Remark 2: One can show here that $n \mathbb{E}(M_n^2) \rightarrow \pi^2/4$.

6.4 Beyond Equivariance: Admissibility

Although equivariance is a rather natural requirement, it is not obvious that it isn't too restrictive. Let us first consider a different estimation paradigm.

Bayesian estimation of θ . Suppose that P_0 is given by a density f_0 on \mathbb{R} , so the distribution of \mathbf{X} is given by the density $f_\theta(\mathbf{x}) = \prod_{i=1}^n f_0(x_i - \theta)$, $\mathbf{x} \in \mathbb{R}^n$. Now imagine that θ itself is a random variable which is chosen by “mother nature” according to a so-called prior distribution with probability density π on \mathbb{R} . That means, for arbitrary Borel sets $C \subset \mathbb{R}$ and $D \subset \mathbb{R}^n$,

$$\mathbb{P}^B(\theta \in C, \mathbf{X} \in D) = \int_C \mathbb{P}_\theta(D) \pi(\theta) d\theta = \int_C \int_D f_\theta(\mathbf{x}) d\mathbf{x} \pi(\theta) d\theta.$$

Here and throughout the sequel, the superscript ‘B’ stands for ‘Bayesian’ and means that θ is considered as a random variable. (We do not distinguish notationally between the random variable θ and an explicit value θ .)

The latter display and Fubini’s theorem show that the joint distribution of (θ, \mathbf{X}) is given by the density $(\theta, \mathbf{x}) \mapsto g(\theta, \mathbf{x}) := f_\theta(\mathbf{x})\pi(\theta)$. That means, for any Borel set $B \subset \mathbb{R} \times \mathbb{R}^n$,

$$\mathbb{P}^B((\theta, \mathbf{X}) \in B) = \int_B g(\theta, \mathbf{x}) d(\theta, \mathbf{x}).$$

Moreover, by Fubini’s theorem,

$$\mathbb{P}^B(\mathbf{X} \in D) = \int_D f^B(\mathbf{x}) d\mathbf{x}$$

with

$$f^B(\mathbf{x}) := \int_{\mathbb{R}} f_\theta(\mathbf{x}) \pi(\theta) d\theta.$$

Hence f^B describes the marginal distribution of \mathbf{X} in the Bayesian framework.

More generally, for any measurable function $h : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\mathbb{E}^B h(\theta, \mathbf{X}) = \int_{\mathbb{R} \times \mathbb{R}^n} h(\theta, \mathbf{x}) g(\theta, \mathbf{x}) d(\theta, \mathbf{x}),$$

provided that the latter integral is well-defined. By Fubini’s theorem this may be rewritten in two ways:

$$\mathbb{E}^B h(\theta, \mathbf{X}) = \int_{\mathbb{R}} \int_{\mathbb{R}^n} h(\theta, \mathbf{x}) f_\theta(\mathbf{x}) d\mathbf{x} \pi(\theta) d\theta = \int_{\mathbb{R}} \mathbb{E}_\theta h(\theta, \mathbf{X}) \pi(\theta) d\theta,$$

and

$$\mathbb{E}^B h(\theta, \mathbf{X}) = \int_{\mathbb{R}^n} \int_{\mathbb{R}} h(\theta, \mathbf{x}) Q_{\mathbf{x}}^B(d\theta) f^B(\mathbf{x}) d\mathbf{x},$$

where $Q_{\mathbf{x}}^B$ is conditional distribution of θ , given $\mathbf{X} = \mathbf{x}$, with density

$$\theta \mapsto \pi(\theta | \mathbf{x}) := \begin{cases} f_\theta(\mathbf{x})\pi(\theta)/f^B(\mathbf{x}) & \text{if } 0 < f^B(\mathbf{x}) < \infty, \\ \pi(\theta) & \text{else.} \end{cases}$$

Within the Bayesian framework, $Q_{\mathbf{x}}^B$ and $\pi(\cdot | \mathbf{x})$ are called the *posterior distribution* and *posterior density*, respectively, of θ , given $\mathbf{X} = \mathbf{x}$.

The *Bayes risk* of any estimator $\hat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$ in this framework is defined as

$$\begin{aligned} R^B(\hat{\theta}) &:= \mathbb{E}^B((\hat{\theta}(\mathbf{X}) - \theta)^2) \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}^n} (\hat{\theta}(\mathbf{x}) - \theta)^2 f_\theta(\mathbf{x}) d\mathbf{x} \pi(\theta) d\theta \\ &= \int_{\mathbb{R}} R(\hat{\theta}, \theta) \pi(\theta) d\theta. \end{aligned}$$

The general theory of conditional expectations implies that the (essentially) unique minimizer of the Bayes risk is given by

$$\widehat{\theta}^B(\mathbf{x}) = \mathbb{E}^B(\theta | \mathbf{X} = \mathbf{x}) = \int_{\mathbb{R}} \theta \pi(\theta | \mathbf{x}) d\theta = \text{mean}(Q_{\mathbf{x}}^B).$$

Furthermore,

$$\begin{aligned} R^B(\widehat{\theta}^B) &= \int_{\mathbb{R} \times \mathbb{R}^n} (\theta - \widehat{\theta}^B(\mathbf{x}))^2 f_{\theta}(\mathbf{x}) \pi(\theta) d(\theta, \mathbf{x}) \\ &= \int_{\mathbb{R}^n} \int_{\mathbb{R}} (\theta - \widehat{\theta}^B(\mathbf{x}))^2 Q_{\mathbf{x}}^B(d\theta) f^B(\mathbf{x}) d\mathbf{x} \\ &= \int \text{Var}(Q_{\mathbf{x}}^B) f^B(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Note the similarity between the Bayes-optimal estimator $\widehat{\theta}^B$ and the optimal equivariant estimator $\widehat{\theta}_*$:

$$\begin{aligned} \widehat{\theta}^B(\mathbf{x}) &= \int_{\mathbb{R}} \theta f_{\theta}(\mathbf{x}) \pi(\theta) d\theta / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) \pi(\theta) d\theta, \\ \widehat{\theta}_*(\mathbf{x}) &= \int_{\mathbb{R}} \theta f_{\theta}(\mathbf{x}) d\theta / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) d\theta. \end{aligned}$$

Hence $\widehat{\theta}_*$ may be interpreted as a Bayesian estimator with prior distribution Lebesgue measure, corresponding to $\pi \equiv 1$. Moreover, suppose that π is the density of $\mathcal{N}(\nu, \tau^2)$ for some $\nu \in \mathbb{R}$ and $\tau > 0$. Then $\pi(\theta)$ is proportional to $e^{-(\theta-\nu)^2/(2\tau^2)}$, whence

$$\widehat{\theta}^B(\mathbf{x}) = \int_{\mathbb{R}} \theta f_{\theta}(\mathbf{x}) e^{-(\theta-\nu)^2/(2\tau^2)} d\theta / \int_{\mathbb{R}} f_{\theta}(\mathbf{x}) e^{-(\theta-\nu)^2/(2\tau^2)} d\theta.$$

Since $(0, 1] \ni e^{-(\theta-\nu)^2/(2\tau^2)} \rightarrow 1$ as $\tau \rightarrow \infty$, it follows from dominated convergence that

$$\widehat{\theta}^B(\mathbf{x}) \rightarrow \widehat{\theta}_*(\mathbf{x}) \quad \text{as } \tau \rightarrow \infty,$$

provided that the integrals in the numerator and denominator of $\widehat{\theta}_*(\mathbf{x})$ are well-defined.

Example 6.20 (Gaussian model and prior). Suppose that π is the density of $\mathcal{N}(0, \tau^2)$ for some $\tau > 0$, and let $P_0 = \mathcal{N}(0, \sigma^2)$ for given $\sigma > 0$. Then

$$\widehat{\theta}^B(\mathbf{x}) = \frac{n}{n + \sigma^2/\tau^2} \bar{x} \quad \text{and} \quad R^B(\widehat{\theta}^B) = \frac{\sigma^2}{n + \sigma^2/\tau^2}.$$

To verify this, recall that

$$f_{\theta}(\mathbf{x}) = f_{\bar{x}}(\mathbf{x}) \exp\left(-\frac{n(\theta - \bar{x})^2}{2\sigma^2}\right).$$

Since $\pi(\theta)$ is proportional to $\exp(-\theta^2/(2\tau^2))$,

$$\begin{aligned} f_{\theta}(\mathbf{x})\pi(\theta) &= C_1(\mathbf{x}) \exp\left(-\frac{\theta^2}{2\tau^2} - \frac{n(\theta - \bar{x})^2}{2\sigma^2}\right) \\ &= C_2(\mathbf{x}) \exp\left(-\frac{\theta^2}{2} \frac{\sigma^2 + n\tau^2}{\sigma^2\tau^2} + \frac{n\bar{x}}{\sigma^2} \theta\right) \\ &= C_3(\mathbf{x}) \exp\left(-\frac{(\theta - \beta\bar{x})^2}{2\gamma^2}\right) \end{aligned}$$

with certain terms $C_1(\mathbf{x}), C_2(\mathbf{x}), C_3(\mathbf{x}) > 0$ and

$$\begin{aligned}\gamma^2 &:= \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} = \frac{\sigma^2}{n + \sigma^2/\tau^2}, \\ \beta &:= \frac{n\gamma^2}{\sigma^2} = \frac{n}{n + \sigma^2/\tau^2}.\end{aligned}$$

In particular, $Q_{\mathbf{x}}^{\text{B}} = \mathcal{N}(\beta\bar{x}, \gamma^2)$, so

$$\widehat{\theta}^{\text{B}}(\mathbf{x}) = \beta\bar{x} \quad \text{and} \quad \text{Var}(Q_{\mathbf{x}}^{\text{B}}) = \gamma^2 = R^{\text{B}}(\widehat{\theta}^{\text{B}}).$$

Admissibility. The use of any estimator $\widehat{\theta}$ is justified if it is admissible in the following sense:

Definition 6.21 (Admissibility). An estimator $\widehat{\theta}$ of θ is called *admissible* if there exists no other estimator $\tilde{\theta}$ such that

$$\begin{aligned}R(\tilde{\theta}, \theta) &\leq R(\widehat{\theta}, \theta) \quad \text{for all } \theta \in \mathbb{R}, \\ R(\tilde{\theta}, \theta_o) &< R(\widehat{\theta}, \theta_o) \quad \text{for some } \theta_o \in \mathbb{R}.\end{aligned}$$

Example 6.22. A rather trivial example of an admissible, but non-equivariant estimator is given by $\widehat{\theta} \equiv \theta_o$ with some fixed value $\theta_o \in \mathbb{R}$, provided that $f_0 > 0$. Here, $R(\widehat{\theta}, \theta) = (\theta_o - \theta)^2$. If $\tilde{\theta}$ would be another estimator with $R(\tilde{\theta}, \cdot) \leq R(\widehat{\theta}, \cdot)$, then $R(\tilde{\theta}, \theta_o) = 0$ is equivalent to $\mathbb{P}_{\theta_o}(\tilde{\theta} \neq \theta_o) = 0$. But for each $\theta \in \mathbb{R}$, the distribution \mathbb{P}_{θ} is absolutely continuous with respect to \mathbb{P}_{θ_o} , so $\mathbb{P}_{\theta}(\tilde{\theta} \neq \theta_o) = 0$, whence $R(\tilde{\theta}, \theta) = (\theta_o - \theta)^2$.

Exercise 6.23. Suppose that $P_0([-1, 1]) = 1$. Show that the trivial estimator $\widehat{\theta} \equiv 0$ is not admissible.

Proposal: Show that if $\mathbf{X} \sim \mathbb{P}_{\theta}$, then $\max(\mathbf{X}) - 1 \leq \theta \leq \min(\mathbf{X}) + 1$ almost surely. Now deduce that $\tilde{\theta}(\mathbf{X}) := (\max(\mathbf{X}) - 1)^+ - (\min(\mathbf{X}) + 1)^-$ outperforms $\widehat{\theta}(\mathbf{X}) = 0$.

The following theorem shows that in case of a centered Gaussian error distribution P_0 , the estimator \bar{X} is indeed admissible.

Theorem 6.24. If $P_0 = \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$, then \bar{X} is an admissible estimator of θ .

Proof of Theorem 6.24. The risk function of \bar{X} is constant σ^2/n . Suppose that $\widehat{\theta}$ is an arbitrary estimator such that $R(\widehat{\theta}, \cdot) \leq \sigma^2/n$ on the whole real line. As shown in Exercises 6.25 and 6.26, it follows from $R(\widehat{\theta}, \cdot) < \infty$ on \mathbb{R} that the risk function $R(\widehat{\theta}, \cdot)$ is continuous. Hence if $R(\widehat{\theta}, \theta_o) < \sigma^2/n$ for some $\theta_o \in \mathbb{R}$, then there exist real numbers $\delta > 0$ and $a < b$ such that $R(\widehat{\theta}, \theta) \leq \sigma^2/n - \delta$ for $\theta \in [a, b]$. Now we evaluate the performance of $\widehat{\theta}$ in a Bayesian framework

with $\theta \sim \mathcal{N}(0, \tau^2)$ for some $\tau > 0$. Here

$$\begin{aligned} R^{\mathbb{B}}(\widehat{\theta}) &= \mathbb{E}^{\mathbb{B}} R(\widehat{\theta}, \theta) \\ &\leq \mathbb{P}^{\mathbb{B}}(\theta \notin [a, b]) \frac{\sigma^2}{n} + \mathbb{P}^{\mathbb{B}}(\theta \in [a, b]) \left(\frac{\sigma^2}{n} - \delta \right) \\ &= \frac{\sigma^2}{n} - \mathbb{P}^{\mathbb{B}}(\theta \in [a, b]) \delta \\ &= \frac{\sigma^2}{n} - \left(\Phi\left(\frac{b}{\tau}\right) - \Phi\left(\frac{a}{\tau}\right) \right) \delta \\ &= \frac{\sigma^2}{n} - \frac{\Phi'(\xi(\tau))(b-a)\delta}{\tau} \end{aligned}$$

for some number $\xi(\tau) \in [a/\tau, b/\tau]$. On the other hand,

$$R^{\mathbb{B}}(\widehat{\theta}) \geq R^{\mathbb{B}}(\widehat{\theta}^{\mathbb{B}}) = \frac{\sigma^2}{n + \sigma^2/\tau^2} \geq \frac{\sigma^2}{n} - \frac{\sigma^4}{n^2\tau^2}$$

by the elementary inequality $1/(1+y) \geq 1-y$ for $y > -1$. These inequalities for $R^{\mathbb{B}}(\widehat{\theta})$ imply that

$$\Phi'(\xi(\tau))(b-a)\delta \leq \frac{\sigma^4}{n^2\tau}$$

for arbitrary $\tau > 0$. But as $\tau \rightarrow \infty$, the left hand side converges to $\Phi'(0)(b-a)\delta > 0$, whereas the right hand side converges to 0. This contradiction shows that $R(\widehat{\theta}, \cdot) \leq \sigma^2/n$ implies that $R(\widehat{\theta}, \cdot) \equiv \sigma^2/n$. \square

Exercise 6.25 (Some basic considerations). Let M be a measure on a measurable space (Ω, \mathcal{A}) , and let $g, h : \Omega \rightarrow \mathbb{R}$ be \mathcal{A} -measurable functions such that for real numbers $a < b$,

$$\int (e^{ag} + e^{bg})|h| dM < \infty.$$

(a) Show that

$$L(t) := \int e^{tg} h dM$$

defines a continuous function $L : [a, b] \rightarrow \mathbb{R}$.

(b) Show that L is continuously differentiable on (a, b) with derivative

$$L'(t) = \int g e^{tg} h dM.$$

(c) Show that L is infinitely often differentiable on (a, b) with k -th derivative

$$L^{(k)}(t) = \int g^k e^{tg} h dM.$$

Exercise 6.26 (Continuity of risk functions in simple Gaussian location families). Consider the simple location family with $P_0 = \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. Let $\widehat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}$ be an estimator of θ such that the risk

$$R(S, \theta) = \mathbb{E}_{\theta}((\widehat{\theta} - \theta)^2)$$

is finite for any $\theta \in \mathbb{R}$. Show that $R(\widehat{\theta}, \cdot)$ is continuous on \mathbb{R} .

Exercise 6.27. Let $Z \sim \mathcal{N}(0, 1)$ and $a \in \mathbb{R}$. Show that $\mathbb{E}(1_{[Z>a]}Z) = \phi(a)$ and $\mathbb{E}(1_{[Z>a]}Z^2) = \Phi(-a) + a\phi(a)$, where ϕ and Φ are the density and distribution function of Z , respectively.

Exercise 6.28. Suppose that the error distribution equals $P_0 = \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. If one assumes that $\theta \geq 0$, a possible estimator would be

$$\bar{X}^+.$$

(a) Determine $R(\bar{X}^+, \theta)$ for arbitrary $\theta \in \mathbb{R}$. (Hint: Exercise 6.27.)

(b) Compare $R(\hat{X}^+, \cdot)$ with $R(\bar{X})$.

Remark: This exercise shows that the estimator \bar{X} of θ is inadmissible in the statistical experiment $(\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (\mathcal{N}(\theta, \sigma^2)^{\otimes n})_{\theta \geq 0})$, because \bar{X}^+ has strictly smaller risk than \bar{X} . Whether or not \bar{X}^+ is admissible itself is a different question.

6.5 Location Functionals and Gross Error Models

Estimators as functionals of (empirical) distributions. Consider a random vector $\mathbf{X} \in \mathbb{R}^n$ with independent components X_i having distribution P . Most estimators $\hat{\theta}(\mathbf{X})$ may be viewed as a functional $S(\hat{P})$ of the empirical distribution

$$\hat{P} := \frac{1}{n} \sum_{i=1}^n \delta_{X_i},$$

i.e.

$$\hat{P}(B) = \frac{\#\{i \leq n : X_i \in B\}}{n} \quad \text{and} \quad \int h d\hat{P} = \frac{1}{n} \sum_{i=1}^n h(X_i)$$

for $B \subset \mathbb{R}$ and $h : \mathbb{R} \rightarrow \mathbb{R}$. For instance

$$\begin{aligned} \bar{X} &= \text{mean}(\hat{P}), \\ \text{median}(X_1, \dots, X_n) &= \text{median}(\hat{P}), \end{aligned}$$

where for arbitrary distributions Q on \mathbb{R} ,

$$\begin{aligned} \text{mean}(Q) &:= \int x Q(dx) \quad \text{provided that} \quad \int |x| Q(dx) < \infty, \\ \text{median}(Q) &:= \frac{\min\{x : Q((-\infty, x]) \geq 0.5\} + \max\{x : Q([x, \infty)) \geq 0.5\}}{2}. \end{aligned}$$

It is well-known that the empirical distribution \hat{P} is a consistent estimator for the underlying distribution P . Precisely,

$$\mathbb{E} \left(\sup_{\text{intervals } B \subset \mathbb{R}} |\hat{P}(B) - P(B)| \right) = O(n^{-1/2})$$

uniformly in P , and for arbitrary measurable functions $h : \mathbb{R} \rightarrow \mathbb{R}$ with $\int |h| dP < \infty$,

$$\mathbb{E} \left| \int h d\hat{P} - \int h dP \right| \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Hence ‘reasonable’ functionals $S(\cdot)$ should satisfy $S(\widehat{P}) \rightarrow_p S(P)$ as $n \rightarrow \infty$, at least if P itself is ‘reasonable’.

In what follows we consider the families

$$\begin{aligned}\mathcal{P} &:= \{\text{probability distributions on } \mathbb{R}\}, \\ \mathcal{P}^r &:= \left\{P \in \mathcal{P} : \int |x|^r P(dx) < \infty\right\}, \quad r \geq 0,\end{aligned}$$

i.e. $\mathcal{P}^0 = \mathcal{P}$.

Definition 6.29 (Equivariant location functional). An equivariant location functional on \mathcal{P}^r is a function $S : \mathcal{P}^r \rightarrow \mathbb{R}$ such that

$$S(P \star \delta_a) = S(P) + a$$

for arbitrary $P \in \mathcal{P}^r$ and $a \in \mathbb{R}$.

Indeed, $\text{mean}(\cdot)$ is an equivariant location functional on \mathcal{P}^1 , and $\text{median}(\cdot)$ is an equivariant location functional on $\mathcal{P}^0 = \mathcal{P}$.

Gross error models. For a given exponent $r \geq 0$ we consider a simple location family

$$(\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (P_\theta^{\otimes n})_{\theta \in \mathbb{R}}) \quad \text{with } P_\theta = P_0 \star \delta_\theta,$$

generated by a given distribution $P_0 \in \mathcal{P}^r$. Now suppose that X_1, X_2, \dots, X_n are independent random variables with distribution P in a ‘contamination neighborhood’ of some distribution in $\{P_\theta : \theta \in \mathbb{R}\}$. Precisely, we assume that for some unknown parameter $\theta \in \mathbb{R}$ and some $\epsilon \in (0, 0.5)$,

$$\begin{aligned}P \in \mathcal{U}_\epsilon^r(\theta) &:= \{(1 - \epsilon)P_\theta + \epsilon Q : Q \in \mathcal{P}^r\} \\ &= \{Q \in \mathcal{P}^r : Q(B) \geq (1 - \epsilon)P_\theta(B) \text{ for any } B \in \text{Borel}(\mathbb{R})\}.\end{aligned}$$

The idea behind this ‘gross error model’ is that each observation X_i stems from P_θ with probability $1 - \epsilon$, but with a (small) probability ϵ it could follow any other distribution $Q \in \mathcal{P}^r$.

For instance, a well-known problem in sociology is that a certain percentage of people give non-sensical answers on questionnaires. In the natural sciences, it may happen that a measurement device fails completely with small probability or that the measured value is recorded with a wrong or missing decimal point which may result in extreme outliers.

If such a model is realistic, for large sample sizes n we should not worry too much about the sampling error $S(\widehat{P}) - S(P)$ but rather about the systematic error $S(P) - \theta$. Note that in case of an equivariant location functional $S : \mathcal{P}^r \rightarrow \mathbb{R}$,

$$\sup_{P \in \mathcal{U}_\epsilon^r(\theta)} |S(P) - \theta| = \sup_{P \in \mathcal{U}_\epsilon^r(0)} |S(P)|$$

for any $\theta \in \mathbb{R}$. For instance, for any $r \geq 1$ and any ‘generator’ $P_0 \in \mathcal{P}^r$ with $\text{mean}(P_0) = 0$,

$$\sup_{P \in \mathcal{U}_\epsilon(0)} |\text{mean}(P)| \geq \sup_{a \in \mathbb{R}} \underbrace{|\text{mean}((1 - \epsilon)P_0 + \epsilon\delta_a)|}_{=\epsilon a} = \infty.$$

Hence the mean is a problematic functional in the presence of gross errors.

The following theorem of Peter J. Huber, a Swiss mathematician and co-founder of the field of “robust statistics”, shows that the median is an optimal equivariant location functional for a broad class of generators P_0 .

Theorem 6.30 (Huber). *For any fixed $r \geq 0$ let $P_0 \in \mathcal{P}^r$ with density f_0 such that f_0 is even on \mathbb{R} and non-increasing on $[0, \infty)$. Then for any equivariant location functional $S : \mathcal{P}^r \rightarrow \mathbb{R}$ and arbitrary $\epsilon \in (0, 0.5)$,*

$$\sup_{P \in \mathcal{U}_\epsilon^r(0)} |S(P)| \geq \sup_{P \in \mathcal{U}_\epsilon^r(0)} |\text{median}(P)| = F_0^{-1}\left(\frac{0.5}{1-\epsilon}\right),$$

where F_0 and F_0^{-1} are the distribution and quantile function, respectively, of P_0 .

Proof of Theorem 6.30. We first show that indeed

$$\sup_{P \in \mathcal{U}_\epsilon(0)^r} |\text{median}(P)| = x_\epsilon := F_0^{-1}\left(\frac{0.5}{1-\epsilon}\right).$$

The assumptions on f_0 imply that with $x_* := \sup\{x > 0 : f_0(x) > 0\} \in (0, \infty]$, the interval $(-x_*, x_*)$ coincides with $\{x \in \mathbb{R} : 0 < F_0(x) < 1\}$, and F_0 is continuous and strictly increasing on $(-x_*, x_*)$. Moreover, $F_0(-x) = 1 - F_0(x)$ for all $x \in \mathbb{R}$. Since $0 < \epsilon < 0.5$, the number x_ϵ lies in $(0, x_*)$. The distribution function F of $P = (1-\epsilon)P_0 + \epsilon Q \in \mathcal{U}_\epsilon^r(P_0)$ is strictly increasing on $(-x_*, x_*)$ as well and satisfies

$$F(-x_\epsilon) \leq (1-\epsilon)F_0(-x_\epsilon) + \epsilon = (1-\epsilon)\left(1 - \frac{0.5}{1-\epsilon}\right) + \epsilon = 0.5$$

with equality if, and only if, $Q((-\infty, -x_\epsilon]) = 1$. On the other hand,

$$F(x_\epsilon) \geq (1-\epsilon)F_0(x_\epsilon) = (1-\epsilon)0.5/(1-\epsilon) = 0.5$$

with equality if, and only if, $Q((-\infty, x_\epsilon]) = 0$. These considerations show that the maximum of $|\text{median}(P)|$ over all $P \in \mathcal{U}_\epsilon^r(0)$ equals x_ϵ .

Now we construct two particular distributions $P^{(1)}, P^{(2)} \in \mathcal{U}_\epsilon^r(0)$ such that

$$P^{(2)} = P^{(1)} \star \delta_{2x_\epsilon}.$$

If this is possible, then for any equivariant location functional $S : \mathcal{R}^r \rightarrow \mathbb{R}$,

$$S(P^{(2)}) - S(P^{(1)}) = 2x_\epsilon.$$

This implies that $S(P^{(1)}) \leq -x_\epsilon$ or $S(P^{(2)}) \geq x_\epsilon$, whence

$$\sup_{P \in \mathcal{U}_\epsilon^r(0)} |S(P)| \geq x_\epsilon.$$

The construction starts from the function $(1-\epsilon)f_0$ and noting that

$$\int_{-\infty}^{x_\epsilon} (1-\epsilon)f_0(x) dx = (1-\epsilon)F_0(x_\epsilon) = 0.5.$$

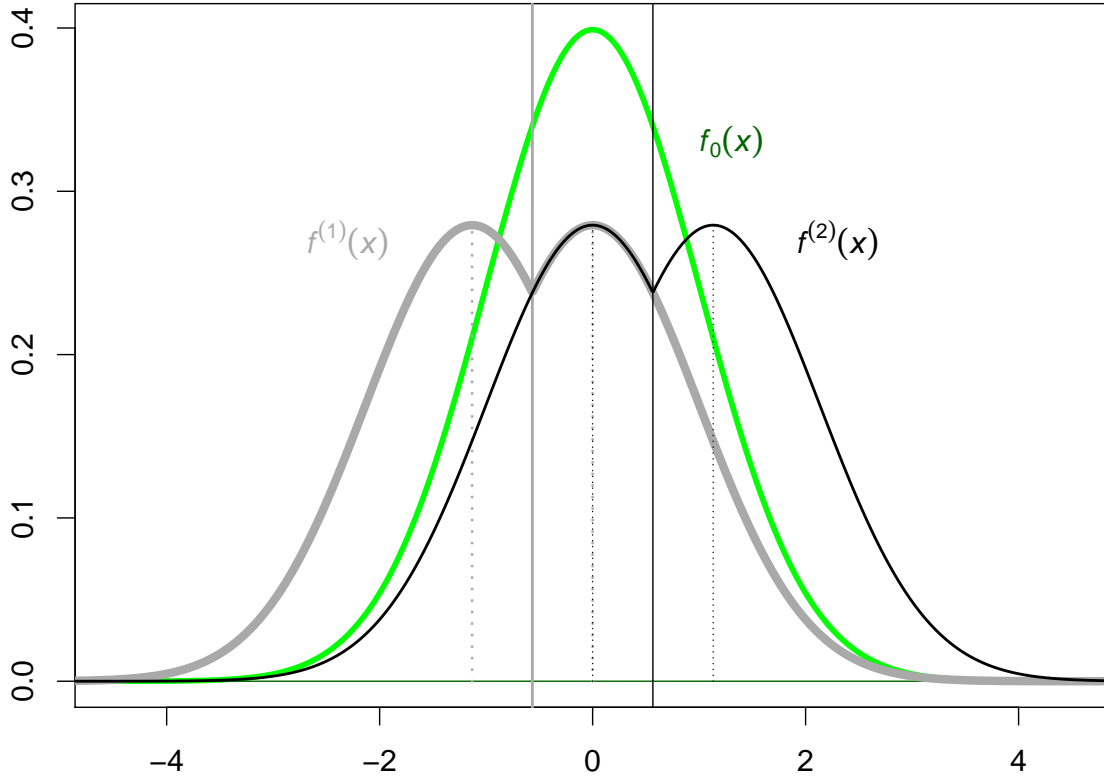


Figure 6.1: Construction of two particular distributions $P^{(1)}, P^{(2)} \in \mathcal{U}_{0.3}(0)$ with densities $f^{(1)}, f^{(2)}$ in case of $P_0 = \mathcal{N}(0, 1)$.

Shifting this function to the right by $2x_\epsilon$ yields the function $(1 - \epsilon)f_0(x - 2x_\epsilon)$, and the assumptions about f_0 imply that

$$(1 - \epsilon)f_0(x) \begin{cases} \geq \\ \leq \end{cases} (1 - \epsilon)f_0(x - 2x_\epsilon) \quad \text{if } x \begin{cases} \leq \\ \geq \end{cases} x_\epsilon,$$

and

$$\int_{x_\epsilon}^{\infty} (1 - \epsilon)f_0(x - 2x_\epsilon) dx = \int_{-x_\epsilon}^{\infty} (1 - \epsilon)f_0(x) dx = 1 - F_0(-x_\epsilon) = 0.5.$$

This shows that

$$f^{(2)} := (1 - \epsilon) \max\{f_0(x), f_0(x - 2x_\epsilon)\}$$

defines a probability density such that the corresponding distribution $P^{(2)}$ belongs to $\mathcal{U}_\epsilon^r(0)$. Instead of shifting $(1 - \epsilon)f_0$ to the right, we could shift it by $2x_\epsilon$ to the left and would obtain the density

$$f^{(1)} := (1 - \epsilon) \max\{f_0(x), f_0(x + 2x_\epsilon)\}$$

of a distribution $P^{(1)} \in \mathcal{U}_\epsilon^r(0)$. But $f^{(2)} = f^{(1)}(\cdot - 2x_\epsilon)$, so $P^{(2)} = P^{(1)} \star \delta_{2x_\epsilon}$, as desired. Figure 6.1 illustrates the construction of $f^{(1)}, f^{(2)}$. \square

Remark. There seems to be no simple location family such that $\text{Median}(x)$ is the corresponding Pitman estimator. On the other hand, if P_0 is the centered Laplace distribution with density

$f_0(x) = (2\sigma)^{-1} \exp(-|x|/\sigma)$, then $\text{Median}(\mathbf{x}) = \hat{\theta}_{\text{ML}}(\mathbf{x})$, and one can show that this estimator is approximately optimal as $n \rightarrow \infty$.

Chapter 7

Statistical Tests

In this chapter we consider a general statistical experiment (also called statistical model)

$$(\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$$

consisting of a sample space (Ω, \mathcal{A}) , a parameter space Θ and given probability distributions \mathbb{P}_θ on (Ω, \mathcal{A}) for arbitrary parameters $\theta \in \Theta$.

Recall that $(\mathbb{P}_\theta)_{\theta \in \Theta}$ describes potential distributions of the observed data $\omega \in \Omega$. Suppose for the moment that the observed data are indeed a realization of a random variable with distribution \mathbb{P}_θ for an unknown *true parameter* θ . Sometimes we conjecture that θ does not belong to a given set $\Theta_o \subset \Theta$. That means, our *working hypothesis* is that $\theta \in \Theta \setminus \Theta_o$, and we would like to falsify the *null hypothesis* that $\theta \in \Theta_o$ based on the observed data. This can be formalized by a measurable function

$$\varphi : \Omega \rightarrow \{0, 1\}.$$

If $\varphi(\omega) = 1$, then we *claim that* $\theta \notin \Theta_o$. In other words we *reject the null hypothesis*. In case of $\varphi(\omega) = 0$ we make *no assertion about* θ .

For theoretical and other reasons it is useful to consider so-called *randomized tests* φ , i.e. measurable mappings

$$\varphi : \Omega \rightarrow [0, 1].$$

The idea is that after observing the data $\omega \in \Omega$, we reject the null hypothesis with (conditional) probability $\varphi(\omega)$. For instance, we could generate an additional random variable $U \sim \text{Unif}[0, 1]$, independently from ω , and reject the null hypothesis if $U \leq \varphi(\omega)$. All in all, by Fubini's theorem, the probability of rejecting the null hypothesis equals

$$\mathbb{E}_\theta(\varphi) = \int \varphi d\mathbb{P}_\theta.$$

This is equal to $\mathbb{P}_\theta(\varphi = 1)$ in case of a $\{0, 1\}$ -valued mapping φ .

Definition 7.1 (Statistical test, power function). A (*statistical*) *test* is a measurable mapping $\varphi : \Omega \rightarrow [0, 1]$. If φ takes only values in $\{0, 1\}$, this can be indicated by saying that φ is a *non-randomized test*. The power function of a test φ is the function

$$\Theta \ni \theta \mapsto \mathbb{E}_\theta(\varphi) = \int \varphi d\mathbb{P}_\theta,$$

and $\mathbb{E}_\theta(\varphi)$ is the *power of φ for parameter θ* .

Note that this definition does not involve any null or working hypothesis. Let us come back to observed data ω coming from \mathbb{P}_θ for an unknown true parameter $\theta \in \Theta$. If we use a test φ to check the null hypothesis that $\theta \in \Theta_o$, there are two possible types of error:

Error of the first kind: The true parameter θ belongs to Θ_o , but we reject the null hypothesis.

Error of the second kind: The true parameter θ does not belong to Θ_o , but we do not reject the null hypothesis.

An error of the first kind happens with probability

$$\begin{cases} \mathbb{E}_\theta(\varphi) & \text{if } \theta \in \Theta_o, \\ 0 & \text{if } \theta \notin \Theta_o, \end{cases}$$

whereas an error of the second kind occurs with probability

$$\begin{cases} 0 & \text{if } \theta \in \Theta_o, \\ 1 - \mathbb{E}_\theta(\varphi) & \text{if } \theta \notin \Theta_o. \end{cases}$$

Traditionally one tries to control the probability of an error of the 1st kind.

Definition 7.2 (Test level). Let $\emptyset \subsetneq \Theta_o \subsetneq \Theta$, and let $\alpha \in (0, 1)$. Suppose that φ is a test such that

$$\mathbb{E}_\theta(\varphi) \leq \alpha \quad \text{for all } \theta \in \Theta_o.$$

Then φ is called a *test of the null hypothesis Θ_o at (test) level α* . A shorter formulation: φ is a *level- α test of Θ_o* .

Example 7.3 (Quality control). The producer of a certain gadget wants to learn something about the unknown probability p that such a device fails in a standardized test of endurance. To this end, he runs an experiment in which n such gadgets are exposed to that endurance test. The outcome of this experiment could be described by a tuple $\omega = (\omega_i)_{i=1}^n$ in $\{0, 1\}^n$, where ω_i specifies whether the i -th gadget fails ($\omega_i = 1$) or not ($\omega_i = 0$). Assuming that the n gadgets perform independently, this leads to the statistical model

$$(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\mathbb{P}_p)_{p \in [0,1]})$$

with \mathbb{P}_p given by

$$\mathbb{P}_p(\{\omega\}) = \prod_{i=1}^n p^{\omega_i} (1-p)^{1-\omega_i} = p^{T(\omega)} (1-p)^{n-T(\omega)}.$$

Here $T(\omega) := \sum_{i=1}^n \omega_i$ is the total number of failures in the experiment.

Alternatively, the producer could focus immediately on the total number of failures in his experiment. Indeed, in a later chapter it will be shown that this reduction is well justified. This leads to the statistical model

$$(\{0, 1, \dots, n\}, \mathcal{P}(\{0, 1, \dots, n\}), (\text{Bin}(n, p))_{p \in [0,1]}).$$

Suppose the producer wants to verify that the unknown probability p is smaller than a given (small) number p_o . Then he should test the null hypothesis $\Theta_o = [p_o, 1]$. If he performs a statistical test of Θ_o at level α , and if that test rejects the null hypothesis, he may claim with confidence $1 - \alpha$ that the unknown parameter p is smaller than p_o .

General goal. Typically we specify a nonempty subset Θ_A of $\Theta \setminus \Theta_o$ and focus on testing the null hypothesis Θ_o versus the alternative hypothesis Θ_A . The goal is to construct a level- α test φ of Θ_o with maximal power $\mathbb{E}_\theta(\varphi)$ for $\theta \in \Theta_A$.

Exercise 7.4 (De-randomisation). Let $\phi : \Omega \rightarrow [0, 1]$ be a statistical test. Show that for any fixed $\beta \in (0, 1)$,

$$\tilde{\phi} := 1_{[\phi \geq \beta]}$$

is a non-randomized test satisfying

$$\frac{\mathbb{E}_\theta(\phi)}{\beta} \geq \mathbb{E}_\theta(\tilde{\phi}) \geq \frac{\mathbb{E}_\theta(\phi) - \beta}{1 - \beta} \quad \text{for all } \theta \in \Theta.$$

7.1 The Neyman–Pearson Lemma

We start with the very simple setting of $\Theta = \{0, 1\}$ and $\Theta_o = \{0\}$.

Theorem 7.5 (Neyman–Pearson). Suppose that \mathbb{P}_0 and \mathbb{P}_1 have densities f_0 and f_1 , respectively, with respect to some measure M on (Ω, \mathcal{A}) . For any $\alpha \in (0, 1)$ there exist constants $k_\alpha \geq 0$ and $\gamma_\alpha \in [0, 1]$ such that

$$\varphi_* := \begin{cases} 1 & \text{if } f_1 > k_\alpha f_0 \\ \gamma_\alpha & \text{if } f_1 = k_\alpha f_0 \\ 0 & \text{if } f_1 < k_\alpha f_0 \end{cases}$$

defines a test φ_* of $\{0\}$ with the following properties:

(i) The test φ_* has exact level α in the sense that

$$\mathbb{E}_0(\varphi_*) = \alpha.$$

(ii) For any level- α test φ of $\{0\}$,

$$\mathbb{E}_1(\varphi) \leq \mathbb{E}_1(\varphi_*).$$

(iii) If φ is a level- α test of $\{0\}$ with $\mathbb{E}_1(\varphi) = \mathbb{E}_1(\varphi_*)$, then

$$M(f_1 > k_\alpha f_0 \text{ and } \varphi < 1) = 0 = M(f_1 < k_\alpha f_0 \text{ and } \varphi > 0).$$

If in addition $k_\alpha > 0$, then $\mathbb{E}_0(\varphi) = \alpha$.

(iv) If $\mathbb{P}_0 \neq \mathbb{P}_1$, then

$$\mathbb{E}_1(\varphi_*) > \alpha.$$

Note that the optimal level- α test φ_* could also be defined in terms of the *likelihood ratio* f_1/f_0 with the conventions that $a/0 := \infty$ for $a > 0$ and $0/0 := 0$:

$$\varphi_* := \begin{cases} 1 & \text{if } f_1/f_0 > k_\alpha, \\ \gamma_\alpha & \text{if } f_1/f_0 = k_\alpha, \\ 0 & \text{if } f_1/f_0 < k_\alpha. \end{cases}$$

Remark 7.6 (Existence and choice of M). The assumption that \mathbb{P}_0 and \mathbb{P}_1 have densities with respect to some measure M on (Ω, \mathcal{A}) is not a real restriction. If we take $M^\circ := \mathbb{P}_0 + \mathbb{P}_1$, then it follows from the theorem of Radon–Nikodym that there exist densities $f_\theta^\circ = d\mathbb{P}_\theta/dM^\circ$ for $\theta = 0, 1$. If M is an arbitrary measure such that a density $f_\theta = d\mathbb{P}_\theta/dM$ exists for $\theta = 0, 1$, then one can easily verify that $f_\theta^\circ = f_\theta/(f_1 + f_2)$ on the set $\{f_1 + f_2 > 0\}$, and $M^\circ(f_1 + f_2 = 0) = 0$. Hence, the resulting optimal test φ_* would be essentially the same, no matter which measure M we start from.

Proof of Theorem 7.5. For the construction of our special test φ_* we consider the auxiliary function $H : [0, \infty) \rightarrow [0, 1]$ given by

$$H(r) = \mathbb{P}_0(f_1 \leq r f_0).$$

One can easily verify that $H(r) = \mathbb{P}_0(f_1/f_0 \leq r)$, where $a/0 := \infty$ for $a > 0$ and $0/0 := 0$. Since $\mathbb{P}_0(f_1/f_0 = \infty) \leq \mathbb{P}_0(f_0 = 0) = 0$, H is a distribution function on $[0, \infty)$. That means, H is nonnegative, nondecreasing and right-continuous with limit $H(+\infty) = 1$. Consequently, the number

$$k_\alpha := \min\{r \geq 0 : H(r) \geq 1 - \alpha\}$$

is well-defined. It has the property that

$$\mathbb{P}_0(f_1 > k_\alpha f_0) = 1 - H(k_\alpha) \leq \alpha \leq 1 - H(k_\alpha -) = \mathbb{P}_0(f_1 \geq k_\alpha f_0).$$

If $\mathbb{P}_0(f_1 = k_\alpha f_0) = 0$, we set $\gamma_\alpha := 1$. Otherwise we define

$$\gamma_\alpha := \frac{\alpha - \mathbb{P}_0(f_1 > k_\alpha f_0)}{\mathbb{P}_0(f_1 \geq k_\alpha f_0) - \mathbb{P}_0(f_1 > k_\alpha f_0)} = \frac{\alpha - \mathbb{P}_0(f_1 > k_\alpha f_0)}{\mathbb{P}_0(f_1 = k_\alpha f_0)} \in [0, 1].$$

In both cases the test $\varphi_* := 1_{[f_1 = k_\alpha f_0]} \gamma_\alpha + 1_{[f_1 > k_\alpha f_0]}$ satisfies

$$\mathbb{E}_0(\varphi_*) = \mathbb{P}_0(f_1 = k_\alpha f_0) \gamma_\alpha + \mathbb{P}_0(f_1 > k_\alpha f_0) = \alpha.$$

This proves property (i).

As to properties (ii-iv), note that for any test φ ,

$$(\varphi - \varphi_*)(f_1 - k_\alpha f_0) \leq 0,$$

because

$$\varphi - \varphi_* = \begin{cases} \varphi - 1 \leq 0 & \text{on } \{f_1 - k_\alpha f_0 > 0\}, \\ \varphi \geq 0 & \text{on } \{f_1 - k_\alpha f_0 < 0\}. \end{cases}$$

Consequently,

$$\begin{aligned}
0 &\geq \int (\varphi - \varphi_*)(f_1 - k_\alpha f_0) dM \\
&= \mathbb{E}_1(\varphi) - \mathbb{E}_1(\varphi_*) - k_\alpha(\mathbb{E}_0(\varphi) - \mathbb{E}_0(\varphi_*)) \\
&= \mathbb{E}_1(\varphi) - \mathbb{E}_1(\varphi_*) - k_\alpha(\mathbb{E}_0(\varphi) - \alpha).
\end{aligned}$$

In other words, for any test φ ,

$$(7.1) \quad \mathbb{E}_1(\varphi) - \mathbb{E}_1(\varphi_*) \leq k_\alpha(\mathbb{E}_0(\varphi) - \alpha)$$

with equality if and only if

$$(7.2) \quad M(f_1 > k_\alpha f_0 \text{ and } \varphi < 1) = 0 = M(f_1 < k_\alpha f_0 \text{ and } \varphi > 0).$$

If φ is a level- α test of $\{0\}$, then the right hand side of (7.1) is non-positive, so $\mathbb{E}_1(\varphi) \leq \mathbb{E}_1(\varphi_*)$. This proves property (ii).

If φ is a level- α test of $\{0\}$ with $\mathbb{E}_1(\varphi) = \mathbb{E}_1(\varphi_*)$, then the right hand side of (7.1) has to be zero, and the first half of property (iii) is just (7.2). Moreover, if $k_\alpha > 0$, then the right hand side of (7.1) being zero means that $\mathbb{E}_0(\varphi) = \alpha$, which proves the second half of property (iii).

Finally, we may compare φ_* with the trivial test $\varphi \equiv \alpha$, so $\mathbb{E}_0(\varphi) = \mathbb{E}_1(\varphi) = \alpha$. Then (7.1) and (7.2) show that $\mathbb{E}_1(\varphi_*) \geq \alpha$ with equality if and only if

$$M(f_1 \neq k_\alpha f_0) = 0.$$

That means, \mathbb{P}_1 has density $k_\alpha f_0$ with respect to M . But then $1 = \mathbb{P}_1(\Omega) = k_\alpha \mathbb{P}_0(\Omega) = k_\alpha$, so $\mathbb{P}_1 = \mathbb{P}_0$. This proves property (iv). \square

Example 7.7. Let $\Omega = (0, \infty)$ and $\mathbb{P}_\theta := \text{Gamma}(a_\theta, b)$ with shape parameters $a_1 > a_0 > 0$ and a common scale parameter $b > 0$. Then the density f_θ of \mathbb{P}_θ with respect to Lebesgue measure on Ω equals $f_\theta(x) = \Gamma(a_\theta)^{-1} b^{-a_\theta} x^{a_\theta-1} e^{-x/b}$, so

$$\frac{f_1}{f_0}(x) = \frac{\Gamma(a_0) b^{a_0-a_1}}{\Gamma(a_1)} x^{a_1-a_0}$$

is strictly increasing in $x > 0$. Hence the optimal level- α test of $\{0\}$ versus $\{1\}$ – in other words: of $\text{Gamma}(a_0, b)$ versus $\text{Gamma}(a_1, b)$ – is given by

$$\varphi_*(x) = 1_{[x \geq k_\alpha]},$$

where k_α is the $(1 - \alpha)$ -quantile of $\text{Gamma}(a_0, b)$.

Exercise 7.8. Let $\mathbb{P}_0 = \mathcal{N}(0, \sigma^2)$ with $\sigma \leq \sqrt{2}$. The corresponding distribution function is $F_0(x) = \Phi(x/\sigma)$. Further let \mathbb{P}_1 be the standard logistic distribution with distribution function

$$F_1(x) = \frac{e^x}{1 + e^x}.$$

Show that the Neyman–Pearson test of \mathbb{P}_0 versus \mathbb{P}_1 (i.e. of $\{0\}$ versus $\{1\}$) at level $\alpha \in (0, 1)$ is given by

$$\varphi_*(x) = 1_{[|x| \geq \sigma \Phi^{-1}(1-\alpha/2)]}.$$

Hint: Show first that

$$\log \frac{f_1(x)}{f_0(x)}$$

is strictly convex and even, where f_0 and f_1 are the density functions of \mathbb{P}_0 and \mathbb{P}_1 , respectively.

Exercise 7.9. Let be $\mathbb{P}_0 = \mathcal{N}(0, 1)$ and $\mathbb{P}_1 = \mathcal{N}(\mu, \sigma^2)$ with $\sigma > 1$.

(a) Show that the Neyman–Pearson test of \mathbb{P}_0 versus \mathbb{P}_1 at level α has the form

$$\varphi_*(x) = 1_{[x \notin (\mu/(1-\sigma^2) \pm \delta_*)]}$$

for some $\delta_* = \delta_*(\mu, \sigma, \alpha) > 0$.

(b) Determine φ_* in the special case of $\mu = 0$.

(c) Show that the special test φ_* in part (b) has the following property:

$$\int \varphi_* d\mathcal{N}(\mu, \sigma^2) \geq \int \varphi_* d\mathcal{N}(0, \sigma^2) \quad \text{for arbitrary } \mu \in \mathbb{R}.$$

Determine the latter power.

7.2 Monotone Density Ratios

In this section we consider a parameter space $\Theta \subset \mathbb{R}$, and we assume that each distribution \mathbb{P}_θ has a density $f_\theta > 0$ with respect to a measure M on (Ω, \mathcal{A}) . Moreover, we assume that there exists a measurable function

$$T : \Omega \rightarrow \mathbb{R}$$

with the following property: For arbitrary $\theta_1, \theta_2 \in \Theta$ with $\theta_1 < \theta_2$, there exists a *non-decreasing function* $g_{\theta_1, \theta_2} : \mathbb{R} \rightarrow (0, \infty)$ such that

$$\frac{f_{\theta_2}}{f_{\theta_1}} = g_{\theta_1, \theta_2}(T).$$

Example 7.10 (Bernoulli experiments). Motivated by Example 7.3, let $\Omega = \{0, 1\}^n$ and $\Theta = (0, 1)$, and let \mathbb{P}_θ describe the joint distribution of n independent random variables with values in $\{0, 1\}$ and expectation θ . That means, with M denoting counting measure on Ω , \mathbb{P}_θ has density

$$f_\theta(\omega) = \prod_{i=1}^n \theta^{\omega_i} (1-\theta)^{1-\omega_i} = \theta^{T(\omega)} (1-\theta)^{n-T(\omega)}$$

with $T(\omega) := \sum_{i=1}^n \omega_i$. Then for $0 < \theta_1 < \theta_2 < 1$ and $\omega \in \Omega$,

$$\begin{aligned} \frac{f_{\theta_2}}{f_{\theta_1}}(\omega) &= \frac{\theta_2^{T(\omega)} (1-\theta_2)^{n-T(\omega)}}{\theta_1^{T(\omega)} (1-\theta_1)^{n-T(\omega)}} \\ &= \frac{(1-\theta_2)^n}{(1-\theta_1)^n} \left(\frac{\theta_2(1-\theta_1)}{(1-\theta_2)\theta_1} \right)^{T(\omega)} \\ &= g_{\theta_1, \theta_2}(T(\omega)) \end{aligned}$$

with and

$$g_{\theta_1, \theta_2}(t) := \frac{(1 - \theta_2)^n (\theta_2(1 - \theta_1))^t}{(1 - \theta_1)^n ((1 - \theta_2)\theta_1)^t}.$$

Note that $g_{\theta_1, \theta_2}(t)$ is strictly increasing in $t \in \mathbb{R}$, because

$$\frac{\theta_2(1 - \theta_1)}{(1 - \theta_2)\theta_1} = \frac{\theta_2}{1 - \theta_2} \bigg/ \frac{\theta_1}{1 - \theta_1} > 1.$$

Example 7.11 (Gaussian location family). Let $\Omega = \mathbb{R}^n$, $\Theta = \mathbb{R}$ and $\mathbb{P}_\theta = \mathcal{N}(\theta, \sigma^2)^{\otimes n}$ for a fixed standard deviation $\sigma > 0$. Recall that the density f_θ of \mathbb{P}_θ with respect to Lebesgue measure on \mathbb{R}^n is given by

$$f_\theta(\mathbf{x}) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\|\mathbf{x} - \theta\|^2}{2\sigma^2}\right) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\|\mathbf{x} - \bar{x}\|^2 + n(\bar{x} - \theta)^2}{2\sigma^2}\right).$$

Thus for $\theta_1 < \theta_2$,

$$\begin{aligned} \frac{f_{\theta_2}(\mathbf{x})}{f_{\theta_1}(\mathbf{x})} &= \exp\left(\frac{n(\bar{x} - \theta_1)^2 - n(\bar{x} - \theta_2)^2}{2\sigma^2}\right) \\ &= \exp\left(\frac{n(\theta_2 - \theta_1)}{\sigma^2} \bar{x} + \frac{n(\theta_1^2 - \theta_2^2)}{2\sigma^2}\right) \\ &= g_{\theta_1, \theta_2}(T(\mathbf{x})), \end{aligned}$$

where $T(\mathbf{x}) := \bar{x}$, and

$$g_{\theta_1, \theta_2}(t) := \exp\left(\frac{n(\theta_2 - \theta_1)}{\sigma^2} t + \frac{n(\theta_1^2 - \theta_2^2)}{2\sigma^2}\right)$$

is strictly increasing in $t \in \mathbb{R}$.

Example 7.12 (Gamma families). As in Example 7.7, let $\Omega = (0, \infty)$, and let Gamma(a, b) be the gamma distribution with shape parameter $a > 0$ and scale parameter $b > 0$. Its density with respect to Lebesgue measure on Ω equals $f_{a,b}(\omega) := f_a(\omega/b)/b$ with $f_a(y) := \Gamma(a)^{-1} \omega^{a-1} e^{-\omega}$. Hence, for parameters $(a_1, b_1), (a_2, b_2) \in (0, \infty) \times (0, \infty)$,

$$\frac{f_{a_2, b_2}(\omega)}{f_{a_1, b_1}(\omega)} = \frac{\Gamma(a_1) b_1^{a_1}}{\Gamma(a_2) b_2^{a_2}} \omega^{a_2 - a_1} \exp((1/b_1 - 1/b_2)\omega),$$

which is strictly increasing in $T(\omega) := \omega$ whenever $a_1 \leq a_2$, $b_1 \leq b_2$ and $(a_1, b_1) \neq (a_2, b_2)$. Consequently, our assumption is satisfied if, for instance, $\Theta = (0, \infty)$ and $\mathbb{P}_\theta = \text{Gamma}(\theta, b)$ for a given $b > 0$ or $\mathbb{P}_\theta = \text{Gamma}(a, \theta)$ for a given $a > 0$.

In statistical models with monotone density ratios as above, there exist optimal tests of null hypotheses of the form

$$\Theta_o = \Theta \cap (-\infty, \theta_o] \quad \text{or} \quad \Theta_o = \Theta \cap [\theta_o, \infty)$$

with arbitrary $\theta_o \in \Theta$.

Theorem 7.13 (Uniformly most powerful (UMP) right-sided tests). *Let $\theta_o \in \Theta$.*

(i) *For any fixed $\alpha \in (0, 1)$ there exist constants $k_\alpha \in \mathbb{R}$ and $\gamma_\alpha \in [0, 1]$ such that the test*

$$\varphi_* := 1_{[T=k_\alpha]} \gamma_\alpha + 1_{[T>k_\alpha]}$$

satisfies

$$\mathbb{E}_{\theta_o}(\varphi_*) = \alpha.$$

(ii) A test φ_* as in part (i) has the following properties:

(ii.1) The power function $\theta \mapsto \mathbb{E}_\theta(\varphi_*)$ is non-decreasing on Θ with values in $(0, 1)$. In particular,

$$\mathbb{E}_\theta(\varphi_*) \leq \alpha \quad \text{for all } \theta \in \Theta \cap (-\infty, \theta_o].$$

(ii.2) For any test φ with $\mathbb{E}_{\theta_o}(\varphi) \leq \alpha$,

$$\mathbb{E}_\theta(\varphi) \leq \mathbb{E}_\theta(\varphi_*) \quad \text{for all } \theta \in \Theta \cap (\theta_o, \infty).$$

(ii.3) For arbitrary parameters $\theta_1 < \theta_2$ with $\mathbb{P}_{\theta_1} \neq \mathbb{P}_{\theta_2}$,

$$\mathbb{E}_{\theta_1}(\varphi_*) < \mathbb{E}_{\theta_2}(\varphi_*).$$

Remark 7.14 (UMP left-sided tests). The previous theorem carries over with obvious modifications to null hypotheses $\Theta_o = \Theta \cap [\theta_o, \infty)$ for some $\theta_o \in \Theta$. Here the optimal level- α test of Θ_o has the form

$$\varphi_* = 1_{[T=k_\alpha]}\gamma_\alpha + 1_{[T < k_\alpha]}$$

with suitable constants $k_\alpha \in \mathbb{R}$ and $\gamma_\alpha \in [0, 1]$.

Proof of Theorem 7.13. The existence of $\gamma_\alpha \in [0, 1]$ and $k_\alpha \in \mathbb{R}$ such that $\varphi_* := 1_{[T=k_\alpha]}\gamma_\alpha + 1_{[T > k_\alpha]}$ satisfies $\mathbb{E}_{\theta_o}(\varphi_*) = \alpha$ can be verified with the same arguments as in the proof of the Neyman–Pearson lemma: We consider the distribution function $H : \mathbb{R} \rightarrow [0, 1]$ with

$$H(r) := \mathbb{P}_{\theta_o}(T \leq r).$$

Then we define

$$k_\alpha := \min\{r \in \mathbb{R} : H(r) \geq 1 - \alpha\},$$

so

$$\mathbb{P}_{\theta_o}(T > k_\alpha) \leq \alpha \leq \mathbb{P}_{\theta_o}(T \geq k_\alpha).$$

In case of $\mathbb{P}_{\theta_o}(T = k_\alpha) = 0$ we set $\gamma_\alpha = 1$, otherwise

$$\gamma_\alpha := \frac{\alpha - \mathbb{P}_{\theta_o}(T > k_\alpha)}{\mathbb{P}_{\theta_o}(T = k_\alpha)} \in (0, 1].$$

Then one can easily verify that the resulting test φ_* has power α at θ_o . This proves part (i).

As to part (ii), we start with a rather general consideration. Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a non-decreasing and bounded function. Then

$$\int h(T) d\mathbb{P}_\theta$$

is a non-decreasing function of $\theta \in \Theta$. For if $\theta_1, \theta_2 \in \Theta$ with $\theta_1 < \theta_2$, then with $g := g_{\theta_1, \theta_2}$,

$$\begin{aligned} \int h(T) d\mathbb{P}_{\theta_2} - \int h(T) d\mathbb{P}_{\theta_1} &= \int h(T)g(T) d\mathbb{P}_{\theta_1} - \int h(T) d\mathbb{P}_{\theta_1} \\ &= \int h(T)(g(T) - 1) d\mathbb{P}_{\theta_1}. \end{aligned}$$

For $h \equiv 1$ we obtain

$$\int (g(T) - 1) d\mathbb{P}_{\theta_1} = 0.$$

Since g is non-decreasing, the latter equation implies that for some $t_o \in \mathbb{R}$,

$$g(t) \begin{cases} \leq 1 & \text{for all } t < t_o, \\ \geq 1 & \text{for all } t > t_o. \end{cases}$$

But then we may conclude that

$$\int h(T)(g(T) - 1) d\mathbb{P}_{\theta_1} = \int (h(T) - h(t_o))(g(T) - 1) d\mathbb{P}_{\theta_1} \geq 0,$$

because the latter integrand is everywhere non-negative.

Note that $\varphi_* = h(T)$ with the nondecreasing function $h : \mathbb{R} \rightarrow [0, 1]$, $h(t) := 1_{[t=k_\alpha]} \gamma_\alpha + 1_{[t>k_\alpha]}$. Consequently, the power function of φ_* is non-decreasing on Θ . Moreover, for arbitrary $\theta \in \Theta$ it follows from $g_{\theta_o, \theta} > 0$ that

$$\int \varphi_* d\mathbb{P}_\theta = \begin{cases} \int \varphi_* g_{\theta_o, \theta} d\mathbb{P}_{\theta_o} > 0, \\ 1 - \int (1 - \varphi_*) g_{\theta_o, \theta} d\mathbb{P}_{\theta_o} < 1, \end{cases}$$

because otherwise $\mathbb{P}_{\theta_o}(\varphi_* > 0) = 0$ or $\mathbb{P}_{\theta_o}(\varphi_* < 1) = 0$, a contradiction to $\mathbb{E}_{\theta_o}(\varphi_*) = \alpha \in (0, 1)$. These considerations prove property (ii.1).

For arbitrary $\theta_1, \theta_2 \in \Theta$ with $\theta_1 < \theta_2$, the function $g := g_{\theta_1, \theta_2}$ is non-decreasing. Hence any test φ satisfies the inequality

$$(\varphi - \varphi_*)(f_{\theta_2} - g(k_\alpha)f_{\theta_1}) = (\varphi - \varphi_*)(g(T) - g(k_\alpha))f_{\theta_1} \leq 0.$$

Consequently

$$\begin{aligned} 0 &\geq \int (\varphi - \varphi_*)(f_{\theta_2} - g(k_\alpha)f_{\theta_1}) dM \\ &= \mathbb{E}_{\theta_2}(\varphi) - \mathbb{E}_{\theta_2}(\varphi_*) - g_{\theta_1, \theta_2}(k_\alpha)(\mathbb{E}_{\theta_1}(\varphi) - \mathbb{E}_{\theta_1}(\varphi_*)), \end{aligned}$$

so

$$(7.3) \quad \mathbb{E}_{\theta_2}(\varphi) - \mathbb{E}_{\theta_2}(\varphi_*) \leq g_{\theta_1, \theta_2}(k_\alpha)(\mathbb{E}_{\theta_1}(\varphi) - \mathbb{E}_{\theta_1}(\varphi_*))$$

with $g_{\theta_1, \theta_2}(k_\alpha) > 0$ by assumption.

In the special case of $\theta_1 = \theta_o$, it follows from (7.3) that $\mathbb{E}_{\theta_2}(\varphi) \leq \mathbb{E}_{\theta_2}(\varphi_*)$ for arbitrary $\theta_2 > \theta_o$ and any test φ satisfying $\mathbb{E}_{\theta_o}(\varphi) \leq \alpha = \mathbb{E}_{\theta_o}(\varphi_*)$. This proves property (ii.2).

As to property (ii.3), (7.3) shows that for arbitrary parameters $\theta_1 < \theta_2$, the test φ_* is an optimal test of the simple null hypothesis $\{\theta_1\}$ versus the simple alternative hypothesis $\{\theta_2\}$ at level $\mathbb{E}_{\theta_1}(\varphi_*) \in (0, 1)$. Thus it follows from the Neyman–Pearson lemma that $\mathbb{E}_{\theta_2}(\varphi_*) > \mathbb{E}_{\theta_1}(\varphi_*)$ unless $\mathbb{P}_{\theta_1} = \mathbb{P}_{\theta_2}$. \square

Example 7.10 (Bernoulli experiments, cont.) Note that the distribution \mathbb{P}_θ^T of T equals the binomial distribution $\text{Bin}(n, \theta)$. Let $b_{n,\theta}$ and $B_{n,\theta}$ denote the weight and distribution function of $\text{Bin}(n, \theta)$, respectively, i.e.

$$b_{n,\theta}(k) := \binom{n}{k} \theta^k (1-\theta)^{n-k},$$

$$B_{n,\theta}(x) := \sum_{k \leq x} b_{n,\theta}(k)$$

for $k, x \in \{0, 1, \dots, n\}$. With the corresponding quantiles

$$B_{n,p}^{-1}(u) := \min\{x : B_{n,p}(x) \geq u\}, \quad u \in (0, 1],$$

for fixed $\theta_o \in (0, 1)$ and $\alpha \in (0, 1)$, the optimal level- α test of $(0, \theta_o]$ versus $(\theta_o, 1)$ is given by

$$\varphi_* = 1_{[T=k_\alpha]} \gamma_\alpha + 1_{[T > k_\alpha]},$$

where

$$k_\alpha = B_{n,\theta_o}^{-1}(1-\alpha) \quad \text{and} \quad \gamma_\alpha = \frac{B_{n,\theta_o}(k_\alpha) - 1 + \alpha}{b_{n,\theta_o}(k_\alpha)}.$$

The power of this test at $\theta \in (0, 1)$ equals

$$\mathbb{E}_\theta(\varphi_*) = 1 - B_{n,\theta}(k_\alpha) + b_{n,\theta}(k_\alpha) \gamma_\alpha.$$

Example 7.11 (Gaussian location family, cont.) Note that in case of $\mathbf{X} \sim \mathcal{N}(\theta, \sigma^2)^{\otimes n}$, the sample mean $\bar{X} = T(\mathbf{X})$ has distribution $\mathcal{N}(\theta, \tau^2)$ with $\tau := \sigma/\sqrt{n}$. Hence,

$$\mathbb{P}_{\theta_o}(T \geq r) = 1 - \Phi\left(\frac{r - \theta_o}{\tau}\right) = \Phi\left(\frac{\theta_o - r}{\tau}\right)$$

equals α if and only if $r = \theta_o - \Phi^{-1}(\alpha)\tau$. Consequently an optimal level- α test of $(-\infty, \theta_o]$ versus (θ_o, ∞) is given by

$$\varphi_*(\mathbf{x}) = 1_{[\bar{x} \geq k_\alpha]} \quad \text{with} \quad k_\alpha = \theta_o - \Phi^{-1}(\alpha)\tau.$$

The power of this test φ_* at an arbitrary parameter θ equals

$$\mathbb{P}_\theta(T \geq k_\alpha) = \Phi\left(\frac{\theta - k_\alpha}{\tau}\right) = \Phi(\Phi^{-1}(\alpha) + (\theta - \theta_o)/\tau).$$

Exercise 7.15. Motivated by the Hardy–Weinberg law in genetics, consider the statistical model $(\mathbb{P}_\theta)_{\theta \in (0,1)}$ with

$$\mathbb{P}_\theta := \text{Mult}(n, \mathbf{p}(\theta)) \quad \text{and} \quad \mathbf{p}(\theta) := ((1-\theta)^2, 2\theta(1-\theta), \theta^2).$$

(a) Show that this model has monotone density ratios for a suitable test statistic $T : \mathbb{N}_0^3 \rightarrow \mathbb{N}_0$.

(b) Determine the distribution \mathbb{P}_θ^T of this test statistic T .

7.3 Stochastic Order, P-Values, Confidence Bounds

In practice, only non-randomized tests, i.e. tests with values in $\{0, 1\}$, are used. Nevertheless the results in Section 7.2 show that certain non-randomized tests based on so-called p-values are essentially optimal. In the present section we describe non-randomized tests, p-values and confidence regions which are valid under a weaker condition on our statistical model $(\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$.

Stochastically ordered distributions. We still assume that Θ is a subset of \mathbb{R} . Further we assume that there exists a measurable function $T : \Omega \rightarrow \mathbb{R}$ such that the corresponding distribution functions $F_\theta : \mathbb{R} \rightarrow [0, 1]$ with

$$F_\theta(t) := \mathbb{P}_\theta(T \leq t)$$

satisfy the following conditions, which are equivalent:

(SO.1) For any fixed $t \in \mathbb{R}$, $F_\theta(t)$ is non-increasing in $\theta \in \Theta$.

(SO.2) For any fixed $t \in \mathbb{R}$, $F_\theta(t-) = \mathbb{P}_\theta(T < t)$ is non-increasing in $\theta \in \mathbb{R}$.

(SO.3) For any fixed $u \in (0, 1)$, $F_\theta^{-1}(u) = \min\{t \in \mathbb{R} : F_\theta(t) \geq u\}$ is non-decreasing in $\theta \in \mathbb{R}$.

(SO.4) For any non-decreasing function $h : \mathbb{R} \rightarrow [0, \infty)$, $\int h(T) d\mathbb{P}_\theta$ is non-decreasing in $\theta \in \mathbb{R}$.

If (SO.1-4) are satisfied, we say that the distribution functions F_θ are stochastically ordered in the sense that $F_{\theta_1} \leq_{\text{st.}} F_{\theta_2}$ whenever $\theta_1 < \theta_2$.

Exercise 7.16. Show that Conditions (SO.1-4) are equivalent.

The proof of Theorem 7.13 (ii) shows that any statistical model with monotone density ratios satisfies condition (SO.4) and hence (SO.1-4). Note also that in Example 7.10 one may extend the parameter space to $\Theta = [0, 1]$, and the stochastic order constraint remains valid.

Tests in terms of p-values. To test whether a hypothetical parameter θ_o is plausible for a particular data set $\omega \in \Omega$, we could compute

- the left-sided p-value

$$F_{\theta_o}(T(\omega)) = \mathbb{P}_{\theta_o}(T \leq T(\omega)),$$

- the right-sided p-value

$$1 - F_{\theta_o}(T(\omega) -) = \mathbb{P}_{\theta_o}(T \geq T(\omega)).$$

The left-sided p-value is non-decreasing in $T(\omega)$, and small values indicate that $T(\omega)$ is “suspiciously small” for the parameter θ_o . The right-sided p-value is non-increasing in $T(\omega)$ with small values indicating that $T(\omega)$ is “suspiciously large” for θ_o .

These p-values lead to tests which are similar to the UMP tests in Section 7.2: Let us fix a test level $\alpha \in (0, 1)$. On the one hand, with the right-sided critical value

$$k_{\alpha, \theta_o}^{(r)} := \min\{t \in \mathbb{R} : F_{\theta_o}(t) \geq 1 - \alpha\}$$

we can conclude that

$$1 - F_{\theta_o}(T-) \leq \alpha \quad \text{if and only if} \quad T \begin{cases} > k_{\alpha, \theta_o}^{(r)} & \text{if } F_{\theta_o}(k_{\alpha, \theta_o}^{(r)}-) < 1 - \alpha, \\ \geq k_{\alpha, \theta_o}^{(r)} & \text{if } F_{\theta_o}(k_{\alpha, \theta_o}^{(r)}-) = 1 - \alpha. \end{cases}$$

Thus by stochastic ordering,

$$\mathbb{P}_\theta(1 - F_{\theta_o}(T-) \leq \alpha) \leq \mathbb{P}_{\theta_o}(1 - F_{\theta_o}(T-) \leq \alpha) \leq \alpha \quad \text{for arbitrary } \theta \leq \theta_o.$$

Consequently,

$$\varphi_{\alpha, \theta_o}^{(r)} := 1_{[1 - F_{\theta_o}(T-) \leq \alpha]} = 1_{[F_{\theta_o}(T-) \geq 1 - \alpha]}$$

defines a level- α test of $\Theta \cap (-\infty, \theta_o]$.

Similarly, with the left-sided critical value

$$k_{\alpha, \theta_o}^{(\ell)} := \max\{t \in \mathbb{R} : F_{\theta_o}(t-) \leq \alpha\}$$

we can conclude that

$$F_{\theta_o}(T) \leq \alpha \quad \text{if and only if} \quad T \begin{cases} < k_{\alpha, \theta_o}^{(\ell)} & \text{if } F_{\theta_o}(k_{\alpha, \theta_o}^{(\ell)}) > \alpha, \\ \leq k_{\alpha, \theta_o}^{(\ell)} & \text{if } F_{\theta_o}(k_{\alpha, \theta_o}^{(\ell)}) = \alpha, \end{cases}$$

whence by stochastic ordering,

$$\mathbb{P}_\theta(F_{\theta_o}(T) \leq \alpha) \leq \mathbb{P}_{\theta_o}(F_{\theta_o}(T) \leq \alpha) \leq \alpha \quad \text{for arbitrary } \theta \geq \theta_o.$$

Consequently,

$$\varphi_{\alpha, \theta_o}^{(\ell)} := 1_{[F_{\theta_o}(T) \leq \alpha]}$$

defines a level- α test of $\Theta \cap [\theta_o, \infty)$.

Confidence bounds. By means of the p-values just constructed, one can also construct confidence regions for the parameter $\theta \in \Theta$:

For given test level $\alpha \in (0, 1)$ and data set $\omega \in \Omega$ let

$$\mathcal{C}_\alpha^{(\ell)}(\omega) := \{\theta \in \Theta : F_\theta(T(\omega)) > \alpha\}.$$

That means, $\mathcal{C}_\alpha^{(\ell)}(\omega)$ is the set of all parameters such that the corresponding left-sided p-value $F_\theta(T(\omega))$ is larger than α . Since $F_\theta(T)$ is non-increasing in $\theta \in \Theta$, the set $\mathcal{C}_\alpha^{(\ell)}(\omega)$ is always an interval $\Theta \cap (-\infty, b_\alpha(\omega))$ or $\Theta \cap (-\infty, b_\alpha(\omega)]$ for some $b_\alpha(\omega) \in [-\infty, \infty]$.

Consequently, $\mathcal{C}_\alpha^{(\ell)}$ and b_α comprise a $(1 - \alpha)$ -confidence interval and an upper $(1 - \alpha)$ -confidence bound in the following sense: For arbitrary $\theta \in \Theta$,

$$\mathbb{P}_\theta(b_\alpha \geq \theta) \geq \mathbb{P}_\theta(\mathcal{C}_\alpha^{(\ell)} \ni \theta) \geq 1 - \alpha.$$

In other words, assuming that an observed data set ω is a realization of a random variable with distribution \mathbb{P}_θ for some unknown true parameter $\theta \in \Theta$, we may claim with confidence $1 - \alpha$ that $\theta \in \mathcal{C}_\alpha^{(\ell)}(\omega)$ and $\theta \leq b_\alpha(\omega)$.

Similarly let

$$\mathcal{C}_\alpha^{(r)}(\omega) := \{\theta \in \Theta : F_\theta(T(\omega) -) < 1 - \alpha\},$$

the set of all parameters such that the corresponding right-sided p-value is larger than α . Since $F_\theta(T(\omega) -)$ is non-increasing in $\theta \in \Theta$, the set $\mathcal{C}_\alpha^{(r)}(\omega)$ is always an interval $\Theta \cap (a_\alpha(\omega), \infty)$ or $\Theta \cap [a_\alpha(\omega), \infty)$ for some $a_\alpha(\omega) \in [-\infty, \infty]$.

Consequently, $\mathcal{C}_\alpha^{(r)}$ and a_α comprise a $(1 - \alpha)$ -confidence interval and a lower $(1 - \alpha)$ -confidence bound in the following sense: For arbitrary $\theta \in \Theta$,

$$\mathbb{P}_\theta(a_\alpha \leq \theta) \geq \mathbb{P}_\theta(\mathcal{C}_{\alpha,r} \ni \theta) \geq 1 - \alpha.$$

That means, assuming that an observed data set ω is a realization of a random variable with distribution \mathbb{P}_θ for some unknown true parameter $\theta \in \Theta$, we may claim with confidence $1 - \alpha$ that $\theta \in \mathcal{C}_\alpha^{(r)}(\omega)$ and $\theta \geq a_\alpha(\omega)$.

Example 7.10 (Bernoulli sequences, cont.) For any $\theta_o \in [0, 1]$ the left- and right-sided p-values are given by

$$B_{n,\theta_o}(T(\omega)) \quad \text{and} \quad 1 - B_{n,\theta_o}(T(\omega) - 1),$$

respectively. The resulting confidence intervals are

$$\mathcal{C}_\alpha^{(\ell)}(\omega) = \begin{cases} [0, b_\alpha(\omega)) & \text{if } T(\omega) < n, \\ [0, 1] & \text{if } T(\omega) = n, \end{cases}$$

$$\mathcal{C}_\alpha^{(r)}(\omega) = \begin{cases} [0, 1] & \text{if } T(\omega) = 0, \\ (a_\alpha(\omega), 1] & \text{if } T(\omega) > 0, \end{cases}$$

where $b_\alpha(\omega)$ is the unique $p \in (0, 1)$ such that $B_{n,p}(T(\omega)) = \alpha$ (if $T(\omega) < n$) while $a_\alpha(\omega)$ is the unique $p \in (0, 1)$ such that $B_{n,p}(T(\omega) - 1) = 1 - \alpha$ (if $T(\omega) > 0$). Here we utilize the fact that for any integer $x \in \{0, 1, \dots, n - 1\}$, the function $p \mapsto B_{n,p}(x)$ is continuous and strictly decreasing with boundary values $B_{n,0}(x) = 1$ and $B_{n,1}(x) = 0$.

Example 7.17 (Gaussian location family, cont.). Since

$$F_\theta(t) = \Phi((t - \theta)/\tau),$$

the left- and right-sided p-values for any given parameter θ_o are given by

$$\Phi((\bar{x} - \theta_o)/\tau) \quad \text{and} \quad \Phi((\theta_o - \bar{x})/\tau),$$

respectively. The resulting confidence intervals are

$$\mathcal{C}_\alpha^{(\ell)}(\mathbf{x}) = (-\infty, b_\alpha(\mathbf{x})) \quad \text{with} \quad b_\alpha(\mathbf{x}) := \bar{x} + \Phi^{-1}(1 - \alpha)\tau,$$

$$\mathcal{C}_\alpha^{(r)}(\mathbf{x}) = (a_\alpha(\mathbf{x}), \infty) \quad \text{with} \quad a_\alpha(\mathbf{x}) := \bar{x} - \Phi^{-1}(1 - \alpha)\tau.$$

Exercise 7.18. Find probability distributions P_0 and P_1 with finite support $\mathcal{X} \subset \mathbb{R}$ and distribution functions F_0 and F_1 , respectively, such that $F_0 \leq_{\text{st.}} F_1$ (i.e. $F_0 \geq F_1$) but $g(x) := P_1(\{x\})/P_0(\{x\})$ is not monotone in $x \in \mathcal{X}$.

Example 7.19 (Capture-recapture). The unknown size N of a population of animals is sometimes estimated with a capture-recapture experiment: At first, a random sample of size n_1 is drawn from the population without replacement, and all animals in this catch are marked and then released. After some time a second sample of size n_2 is drawn without replacement, and one determines the number X of marked animals in this second catch. That means, X is the number of animals which were caught twice. Ideally, X is a random variable with distribution $\text{Hyp}(N, n_1, n_2)$. This leads to the statistical experiment

$$(\mathcal{X}, \mathcal{P}(\mathcal{X}), (\text{Hyp}(N, n_1, n_2))_{N \geq \max(n_1, n_2)}),$$

where $\mathcal{X} := \{0, 1, \dots, \min(n_1, n_2)\}$ and

$$\text{Hyp}(N, n_1, n_2)(\{x\}) = \binom{n_1}{x} \binom{N - n_1}{n_2 - x} / \binom{N}{n_2} = \binom{n_2}{x} \binom{N - n_2}{n_1 - x} / \binom{N}{n_1}$$

for $x \in \mathcal{X}$ with the convention that $\binom{k}{\ell} := 0$ if $\ell > k$. Possible point estimators for N are given by

$$\widehat{N}(x) := \frac{n_1 n_2}{x} \quad \text{or} \quad \widehat{N}(x) := \frac{(n_1 + 1)(n_2 + 1)}{x + 1}.$$

It follows from Exercise 7.20 below that the distributions $\text{Hyp}(N, n_1, n_2)$ are stochastically *decreasing* in N . That means, if F_N denotes the distribution function of $\text{Hyp}(N, n_1, n_2)$, then

$$F_N(x) \text{ is non-decreasing in } N \geq \max(n_1, n_2)$$

for any $x \in \mathcal{X}$, and

$$\lim_{N \rightarrow \infty} F_N(0) = 1.$$

Consequently, a lower $(1 - \alpha)$ -confidence bound for N is given by

$$a_\alpha(x) := \min\{N \geq \max(n_1, n_2) : F_N(x) > \alpha\}$$

while an upper bound is given by

$$b_\alpha(x) = \begin{cases} \infty & \text{if } x = 0, \\ \max\{N \geq \max(n_1, n_2) : F_N(x - 1) < 1 - \alpha\} & \text{if } x > 0. \end{cases}$$

Exercise 7.20. For integers $n_1, n_2 \geq 1$ and $N \geq \max(n_1, n_2)$, let F_N be the distribution function of the hypergeometric distribution $\text{Hyp}(N, n_1, n_2)$.

(a) Show that $F_N(x) \leq F_{N+1}(x)$ for arbitrary $x \in \{0, 1, \dots, \min(n_1, n_2)\}$. Proposal: Think of an urn with n_1 black, $N - n_1$ white and one red ball from which you draw $n_2 + 1$ balls one by one without replacement.

(b) Show that

$$\lim_{N \rightarrow \infty} F_N(0) = 1.$$

Exercise 7.21. Let X be a random variable with distribution $\text{Hyp}(N, n_1, n_2)$ with given parameters $n_1, n_2 \in \mathbb{N}$ and an unknown parameter $N \geq \max(n_1, n_2)$. Determine $\mathbb{E}_N(\widehat{N})$ for the point estimator $\widehat{N} := (n_1 + 1)(n_2 + 1)/(X + 1)$.

7.4 The Generalized Neyman–Pearson Lemma

Our goal is to construct optimal tests of null hypotheses Θ_o such that $\#\Theta_o > 1$. In the setting of monotone density ratios, we solved this problem for $\Theta_o = \Theta \cap (-\infty, \theta_o]$ or $\Theta_o = \Theta \cap [\theta_o, \infty)$ with a “least favourable parameter” $\theta_o \in \Theta$. But what about null hypotheses without such a unique least favourable parameter? To deal with such settings, we start with a rather general consideration.

Theorem 7.22 (Generalized Neyman–Pearson lemma). *Let \mathcal{T} be the set of all statistical tests $\varphi : \Omega \rightarrow [0, 1]$. Let M be a σ -finite measure on (Ω, \mathcal{A}) , and let $f_1, \dots, f_m, f_{m+1} \in \mathcal{L}^1(M)$ for some integer $m \geq 1$. Further let $\alpha \in \mathbb{R}^m$ and define*

$$\mathcal{T}(\alpha) := \left\{ \varphi \in \mathcal{T} : \int \varphi \mathbf{f} dM = \alpha \right\}$$

with $\mathbf{f} = (f_j)_{j=1}^m : \Omega \rightarrow \mathbb{R}^m$.

(i) *If $\mathcal{T}(\alpha) \neq \emptyset$, then there exists a test $\varphi_* \in \mathcal{T}(\alpha)$ such that*

$$\int \varphi_* f_{m+1} dM \geq \int \varphi f_{m+1} dM \quad \text{for all } \varphi \in \mathcal{T}(\alpha).$$

(ii) *Suppose that φ_* is a test in $\mathcal{T}(\alpha)$ such that*

$$\varphi_*(\omega) = \begin{cases} 1 & \text{if } f_{m+1}(\omega) > \mathbf{k}_\alpha^\top \mathbf{f}(\omega) \\ 0 & \text{if } f_{m+1}(\omega) < \mathbf{k}_\alpha^\top \mathbf{f}(\omega) \end{cases}$$

for a certain $\mathbf{k}_\alpha = (k_{\alpha,j})_{j=1}^m \in \mathbb{R}^m$. Then φ_* has the optimality property in part (i). More generally,

$$\int \varphi_* f_{m+1} dM \geq \int \varphi f_{m+1} dM$$

for arbitrary tests $\varphi \in \mathcal{T}$ such that for $1 \leq j \leq m$,

$$\int \varphi f_j dM \begin{cases} \geq \alpha_j & \text{if } k_{\alpha,j} < 0, \\ \leq \alpha_j & \text{if } k_{\alpha,j} > 0. \end{cases}$$

(iii) *Suppose that α is an interior point of the set*

$$\left\{ \int \varphi \mathbf{f} dM : \varphi \in \mathcal{T} \right\} \subset \mathbb{R}^m.$$

Then there exists a test φ_* as described in part (ii).

Proof of Theorem 7.22. In what follows let

$$\begin{aligned} \mathcal{K}_m &:= \left\{ \int \varphi \mathbf{f} dM : \varphi \in \mathcal{T} \right\} \subset \mathbb{R}^m, \\ \mathcal{K}_{m+1} &:= \left\{ \left(\int \varphi \mathbf{f} dM, \int \varphi f_{m+1} dM \right) : \varphi \in \mathcal{T} \right\} \subset \mathbb{R}^m \times \mathbb{R}. \end{aligned}$$

The set \mathcal{K}_{m+1} is a compact and convex subset of $\mathbb{R}^m \times \mathbb{R}$. This can be verified in two different ways:

With $\mathcal{F} := \mathcal{L}^1(M)$, the set

$$\mathcal{K} := \left\{ \left(\int \varphi f dM \right)_{f \in \mathcal{F}} : \varphi \in \mathcal{T} \right\}$$

is a compact and convex subset of $\mathbb{R}^{\mathcal{F}}$, equipped with the product topology, see Theorem A.12 in Appendix A. But \mathcal{K}_{m+1} is the image of \mathcal{K} under the linear and continuous mapping

$$\mathbb{R}^{\mathcal{F}} \ni (x_f)_{f \in \mathcal{F}} \mapsto ((x_{f_j})_{j=1}^m, x_{f_{m+1}}) \in \mathbb{R}^m \times \mathbb{R},$$

whence it is a compact and convex subset of $\mathbb{R}^m \times \mathbb{R}$.

Alternatively, one can verify directly that \mathcal{K}_{m+1} is a convex and bounded subset of $\mathbb{R}^m \times \mathbb{R}$. But Theorem A.14 implies that it is closed, whence it is compact.

Proof of part (i): Since \mathcal{K}_{m+1} is compact and convex, its intersection with the set

$$\{\boldsymbol{\alpha}\} \times \mathbb{R}$$

is empty or of the form

$$\{\boldsymbol{\alpha}\} \times [a, b]$$

with real numbers $a \leq b$. In the latter case there exists a test $\varphi_* \in \mathcal{T}(\boldsymbol{\alpha})$ such that

$$\int \varphi_* f_{m+1} dM = b = \max_{\varphi \in \mathcal{T}(\boldsymbol{\alpha})} \int \varphi f_{m+1} dM.$$

Proof of part (ii): Let $\varphi_* \in \mathcal{T}(\boldsymbol{\alpha})$ have the specified special form. Then for any other test φ ,

$$(\varphi_* - \varphi)(f_{m+1} - \mathbf{k}_{\boldsymbol{\alpha}}^{\top} \mathbf{f}) \geq 0,$$

whence

$$\begin{aligned} \int \varphi_* f_{m+1} dM - \int \varphi f_{m+1} dM &= \int (\varphi_* - \varphi)(f_{m+1} - \mathbf{k}_{\boldsymbol{\alpha}}^{\top} \mathbf{f}) dM \\ &\quad + \sum_{j=1}^m k_{\boldsymbol{\alpha},j} \left(\alpha_j - \int \varphi f_j dM \right) \\ &\geq \sum_{j=1}^m k_{\boldsymbol{\alpha},j} \left(\alpha_j - \int \varphi f_j dM \right). \end{aligned}$$

The right hand side equals 0, if $\varphi \in \mathcal{T}(\boldsymbol{\alpha})$, so φ_* maximizes $\int \varphi f_{m+1} dM$ over all tests $\varphi \in \mathcal{T}(\boldsymbol{\alpha})$. More generally, the right hand side is non-negative for all tests φ such that for $1 \leq j \leq m$,

$$\int \varphi f_j dM \begin{cases} \leq \alpha_j & \text{if } k_{\boldsymbol{\alpha},j} > 0, \\ \geq \alpha_j & \text{if } k_{\boldsymbol{\alpha},j} < 0. \end{cases}$$

Proof of part (iii): Let

$$\mathcal{C} := \{\boldsymbol{\alpha}\} \times (b, \infty)$$

with

$$b := \max_{\varphi \in \mathcal{T}(\boldsymbol{\alpha})} \int \varphi f_{m+1} dM.$$

Then \mathcal{K}_{m+1} and \mathcal{C} are disjoint convex subsets of $\mathbb{R}^m \times \mathbb{R}$. Consequently they may be separated weakly by a hyperplane. That means, there exists a nonzero vector $(\mathbf{k}, u) \in \mathbb{R}^m \times \mathbb{R}$ such that

$$\langle (\mathbf{x}, y), (\mathbf{k}, u) \rangle \leq \langle (\boldsymbol{\alpha}, z), (\mathbf{k}, u) \rangle \quad \text{for arbitrary } (\mathbf{x}, y) \in \mathcal{K}_{m+1} \text{ and } z > b$$

with $\langle \cdot, \cdot \rangle$ denoting the standard inner product on $\mathbb{R}^m \times \mathbb{R}$. Thus

$$\mathbf{k}^\top \mathbf{x} + uy \leq \mathbf{k}^\top \boldsymbol{\alpha} + uz \quad \text{for arbitrary } (\mathbf{x}, y) \in \mathcal{K}_{m+1} \text{ and } z > b.$$

Fixing one point $(\mathbf{x}, y) \in \mathcal{K}_{m+1}$ and letting $z \rightarrow \infty$ shows that $u \geq 0$. In case of $u = 0$, we would have $\mathbf{k} \neq \mathbf{0}$ and

$$\mathbf{k}^\top \mathbf{x} \leq \mathbf{k}^\top \boldsymbol{\alpha} \quad \text{for arbitrary } \mathbf{x} \in \mathcal{K}_m.$$

But then $\boldsymbol{\alpha}$ would be a boundary point of \mathcal{K}_m , rather than an interior point. Consequently, $u > 0$, and we may assume without loss of generality that $u = 1$. Consequently, for arbitrary tests $\varphi \in \mathcal{T}$,

$$\mathbf{k}^\top \int \varphi \mathbf{f} dM + \int \varphi f_{m+1} dM \leq \mathbf{k}^\top \boldsymbol{\alpha} + b.$$

If $\varphi_* \in \mathcal{T}(\boldsymbol{\alpha})$ with $\int \varphi_* f_{m+1} dM = b$, then with $\mathbf{k}_\alpha := -\mathbf{k}$, we may rewrite the previous inequality as

$$\int (\varphi_* - \varphi)(f_{m+1} - \mathbf{k}_\alpha^\top \mathbf{f}) dM \geq 0$$

for arbitrary tests φ . Applying this inequality to the special test

$$\varphi_{**} := \begin{cases} 1 & \text{if } f_{m+1} > \mathbf{k}_\alpha^\top \mathbf{f} \\ 0 & \text{if } f_{m+1} < \mathbf{k}_\alpha^\top \mathbf{f} \\ \varphi_* & \text{else} \end{cases}$$

shows that

$$M(\varphi_* \neq \varphi_{**}) = 0.$$

Hence we may replace φ_* with φ_{**} . □

7.5 Tests of Two-Sided Hypotheses

7.5.1 One-parameter exponential families (with natural parametrization)

In what follows we apply the generalized Neyman–Pearson lemma to a particular type of statistical model $(\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$: Let Θ be a real interval, and suppose that for some σ -finite measure M on (Ω, \mathcal{A}) ,

$$f_\theta(\omega) = \frac{d\mathbb{P}_\theta}{dM}(\omega) = C(\theta)h(\omega) \exp(\theta T(\omega))$$

for given measurable functions $h : \Omega \rightarrow [0, \infty)$, $T : \Omega \rightarrow \mathbb{R}$ and the normalisation constant

$$C(\theta) = \left(\int h e^{\theta T} dM \right)^{-1}.$$

We also assume that $M(h > 0 \text{ and } T \neq c) > 0$ for any real constant c , so $\mathbb{P}_{\theta_1} \neq \mathbb{P}_{\theta_2}$ for arbitrary different $\theta_1, \theta_2 \in \Theta$.

Example 7.10 (Bernoulli sequences, cont.) This statistical model is an exponential family with $\Omega = \{0, 1\}^n$, $M =$ counting measure on Ω , and

$$\begin{aligned}\Theta &:= \mathbb{R}, \\ \theta(p) &:= \log(p/(1-p)), \\ h(\omega) &:= 1, \\ T(\omega) &:= \sum_{i=1}^n \omega_i, \\ C(\theta) &:= (1 + e^\theta)^{-n}, \quad \text{i.e. } C(\theta(p)) = (1-p)^n.\end{aligned}$$

Example 7.23 (Poisson distributions). The family of Poisson distributions $\text{Poiss}(\lambda)$, $\lambda > 0$, is an exponential family with $\Omega = \mathbb{N}_0$, $M =$ counting measure on Ω , and

$$\begin{aligned}\Theta &:= \mathbb{R}, \\ \theta(\lambda) &:= \log(\lambda), \\ h(\omega) &:= (\omega!)^{-1}, \\ T(\omega) &:= \omega, \\ C(\theta) &:= \exp(-e^\theta), \quad \text{i.e. } C(\theta(\lambda)) = e^{-\lambda}.\end{aligned}$$

Example 7.11 (Gaussian location family, cont.) The family of distributions $\mathcal{N}(\mu, \sigma^2)^{\otimes n}$, $\mu \in \mathbb{R}$, is an exponential family with $\Omega = \mathbb{R}^n$, $M =$ Lebesgue measure on \mathbb{R}^n , and

$$\begin{aligned}\Theta &:= \mathbb{R}, \\ \theta(\mu) &:= n\mu/\sigma^2, \\ h(\mathbf{x}) &:= (2\pi\sigma^2)^{-1/2} \exp(-\|\mathbf{x}\|^2/(2\sigma^2)), \\ T(\mathbf{x}) &:= \bar{x}, \\ C(\theta) &:= \exp(-\sigma^2\theta^2/(2n)).\end{aligned}$$

Alternatively, one could choose, for instance,

$$\begin{aligned}\theta(\mu) &:= \mu, \\ T(\mathbf{x}) &:= n\bar{x}/\sigma^2, \\ C(\theta) &:= \exp(-n\theta^2/(2\sigma^2)),\end{aligned}$$

or

$$\begin{aligned}\theta(\mu) &:= \sqrt{n}\mu/\sigma, \\ T(\mathbf{x}) &:= \sqrt{n}\bar{x}/\sigma, \\ C(\theta) &:= \exp(-\theta^2/2).\end{aligned}$$

An advantage of the latter parametrization is that $\mathbb{P}_\theta^T = \mathcal{N}(\theta, 1)$.

7.5.2 Two-sided hypotheses, version 1

Now we consider the problem of testing

$$\Theta_o := \{\theta \in \Theta : \theta \leq \theta_1 \text{ or } \theta \geq \theta_2\} \quad \text{versus} \quad \Theta_A := (\theta_1, \theta_2)$$

with given parameters $\theta_1, \theta_2 \in \Theta$ such that $\theta_1 < \theta_2$.

Theorem 7.24. (i) For any fixed $\alpha \in (0, 1)$ there exist real constants $c_1 \leq c_2$ and $\gamma_1, \gamma_2 \in [0, 1]$ (with $\gamma_2 = 0$ in case of $c_1 = c_2$) such that

$$\varphi_* := 1_{[T=c_1]}\gamma_1 + 1_{[T=c_2]}\gamma_2 + 1_{[c_1 < T < c_2]}$$

is a test satisfying

$$\mathbb{E}_{\theta_1}(\varphi_*) = \mathbb{E}_{\theta_2}(\varphi_*) = \alpha.$$

(ii) A test φ_* as in part (i) has the following properties:

(ii.1) For any level- α test φ of $\{\theta_1, \theta_2\}$,

$$\mathbb{E}_{\theta}(\varphi_*) \geq \mathbb{E}_{\theta}(\varphi) \quad \text{for all } \theta \in (\theta_1, \theta_2).$$

(ii.2) For arbitrary tests φ such that $\mathbb{E}_{\theta_1}(\varphi) = \mathbb{E}_{\theta_2}(\varphi) = \alpha$,

$$\mathbb{E}_{\theta}(\varphi_*) \leq \mathbb{E}_{\theta}(\varphi) \quad \text{for all } \theta \in \Theta_o.$$

In particular, φ_* is a level- α test of Θ_o , i.e.

$$\mathbb{E}_{\theta}(\varphi_*) \leq \alpha \quad \text{for all } \theta \in \Theta_o.$$

For the proof of this and later theorems, we need an elementary result about weighted sums of two exponential functions.

Lemma 7.25. For real numbers $c_1 \leq c_2$ and $d_1 < d_2$ there exists a unique vector $\mathbf{b} \in \mathbb{R}^2$ such that the function $A : \mathbb{R} \rightarrow \mathbb{R}$,

$$A(t) := \sum_{j=1}^2 b_j e^{d_j t}$$

satisfies

$$\begin{cases} A(c_1) = A(c_2) = 1, \\ A'(c_1) = 0 \end{cases} \quad \text{if } c_1 = c_2.$$

If $d_1 < 0 < d_2$, then $b_1, b_2 > 0$ and

$$A \begin{cases} < 1 & \text{on } (c_1, c_2), \\ > 1 & \text{on } \mathbb{R} \setminus [c_1, c_2]. \end{cases}$$

If $0 < d_1 < d_2$, then $b_1 > 0 > b_2$, whereas $d_1 < d_2 < 0$ implies that $b_1 < 0 < b_2$. In both cases,

$$A \begin{cases} > 1 & \text{on } (c_1, c_2), \\ < 1 & \text{on } \mathbb{R} \setminus [c_1, c_2]. \end{cases}$$

Proof of Lemma 7.25. Suppose first that $c_1 < c_2$. The condition $A(c_1) = A(c_2) = 1$ is equivalent to

$$\mathbf{A}\mathbf{b} = (1, 1)^\top$$

with the matrix

$$\mathbf{A} := \begin{bmatrix} e^{d_1 c_1} & e^{d_2 c_1} \\ e^{d_1 c_2} & e^{d_2 c_2} \end{bmatrix}.$$

Note that

$$\det(\mathbf{A}) = e^{d_1 c_1 + d_2 c_2} - e^{d_2 c_1 + d_1 c_2} = e^{d_2 c_1 + d_1 c_2} (e^{(d_2 - d_1)(c_2 - c_1)} - 1) > 0.$$

Hence the equation $\mathbf{A}\mathbf{b} = (1, 1)^\top$ has the unique solution

$$\mathbf{b} = \det(\mathbf{A})^{-1} \begin{bmatrix} e^{d_2 c_2} & -e^{d_2 c_1} \\ -e^{d_1 c_2} & e^{d_1 c_1} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \det(\mathbf{A})^{-1} \begin{bmatrix} e^{d_2 c_1} (e^{d_2(c_2 - c_1)} - 1) \\ e^{d_1 c_2} (e^{-d_1(c_2 - c_1)} - 1) \end{bmatrix},$$

and the stated inequalities for b_1, b_2 are clearly satisfied.

In case of $c_1 = c_2 = c$, the equations $A(c) = 1$ and $A'(c) = 0$ are equivalent to

$$\mathbf{A}\mathbf{b} = (1, 0)^\top$$

with the matrix

$$\mathbf{A} := \begin{bmatrix} e^{d_1 c} & e^{d_2 c} \\ d_1 e^{d_1 c} & d_2 e^{d_2 c} \end{bmatrix}.$$

Again, $\det(\mathbf{A}) = (d_2 - d_1)e^{(d_1 + d_2)c} > 0$, so the equation $\mathbf{A}\mathbf{b} = (1, 0)^\top$ has the unique solution

$$\mathbf{b} = \det(\mathbf{A})^{-1} \begin{bmatrix} d_2 e^{d_2 c} & -e^{d_2 c} \\ -d_1 e^{d_1 c} & e^{d_1 c} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \det(\mathbf{A})^{-1} \begin{bmatrix} e^{d_2 c} d_2 \\ -e^{d_1 c} d_1 \end{bmatrix},$$

and the stated inequalities for b_1, b_2 are clearly satisfied.

In case of $b_1, b_2 > 0$ the function A is strictly convex, whence $A < 1$ on (c_1, c_2) and $A > 1$ on $\mathbb{R} \setminus [c_1, c_2]$.

In case of $0 < d_1 < d_2$ and $b_1 > 0 > b_2$,

$$A'(t) = b_1 d_1 e^{d_1 t} + b_2 d_2 e^{d_2 t} = |b_2| d_2 e^{d_1 t} \left(\frac{b_1 d_1}{|b_2| d_2} - e^{(d_2 - d_1)t} \right) \begin{cases} > 0 & \text{if } t < t_o \\ < 0 & \text{if } t > t_o \end{cases}$$

for some $t_o \in \mathbb{R}$. If $c_1 < c_2$, it follows from $A(c_1) = A(c_2) = 1$ that $t_o \in (c_1, c_2)$, whence $\{A > 1\} = (c_1, c_2)$ and $\{A < 1\} = \mathbb{R} \setminus [c_1, c_2]$. If $c_1 = c_2$, it follows from $A'(c_1) = 0$ that $t_o = c_1$, whence $\{A > 1\} = \emptyset$ and $\{A < 1\} = \mathbb{R} \setminus \{c_1\}$.

Analogous considerations apply in case of $d_1 < d_2 < 0$ and $b_1 < 0 < b_2$. □

Proof of Theorem 7.24. Since $\mathbb{P}_\theta(h = 0) = 0$ for all $\theta \in \Theta$, we may replace Ω with $\{h > 0\}$ and $M(d\omega)$ with $h(\omega)M(d\omega)$, so $h \equiv 1$.

Proof of part (i). We fix an arbitrary parameter $\theta_3 \in (\theta_1, \theta_2)$ and construct an optimal level- α test φ_* of $\{\theta_1, \theta_2\}$ versus $\{\theta_3\}$ by means of the generalized Neyman–Pearson lemma. To this end we consider the point $\alpha := (\alpha, \alpha)^\top$ and the set

$$\mathcal{K}_2 = \left\{ \int \varphi \mathbf{f} dM : \varphi \in \mathcal{T} \right\}$$

with $\mathbf{f} := (f_{\theta_1}, f_{\theta_2})^\top : \Omega \rightarrow \mathbb{R}^2$. Let $\varphi_0 \equiv 0$, let φ_1 be an optimal level- α test of $\{\theta_1\}$ versus $\{\theta_2\}$, and let φ_2 be an optimal level- α test of $\{\theta_2\}$ versus $\{\theta_1\}$. Then

$$\begin{aligned} \int \varphi_0 \mathbf{f} dM &= (0, 0)^\top, \\ \int \varphi_1 \mathbf{f} dM &= (\alpha, \alpha_1)^\top \quad \text{for some } \alpha_1 > \alpha, \\ \int \varphi_2 \mathbf{f} dM &= (\alpha_2, \alpha)^\top \quad \text{for some } \alpha_2 > \alpha. \end{aligned}$$

This implies that α is an interior point of the set \mathcal{K}_2 , see also Exercise 7.26 below. Consequently, there exists a test φ_* such that $\mathbb{E}_{\theta_1}(\varphi_*) = \mathbb{E}_{\theta_2}(\varphi_*) = \alpha$ with

$$\varphi_* = \begin{cases} 1 & \text{if } f_{\theta_3} > k_1 f_{\theta_1} + k_2 f_{\theta_2} \\ 0 & \text{if } f_{\theta_3} < k_1 f_{\theta_1} + k_2 f_{\theta_2} \end{cases}$$

for certain constants $k_1, k_2 \in \mathbb{R}$. With $b_j := k_j C(\theta_j)/C(\theta_3)$ we may also write

$$\varphi_* = \begin{cases} 1 & \text{if } A(T) < 1 \\ 0 & \text{if } A(T) > 1 \end{cases}$$

with

$$A(t) := b_1 e^{(\theta_1 - \theta_3)t} + b_2 e^{(\theta_2 - \theta_3)t}.$$

Since $\mathbb{E}_{\theta_1}(\varphi_*), \mathbb{E}_{\theta_2}(\varphi_*) < 1$, we may conclude that $\max(b_1, b_2) > 0$. But then the function A is strictly monotone or strictly convex. Hence the set $\{A = 1\}$ has at most two elements, and we may replace φ_* with

$$\bar{\varphi}_* := \begin{cases} 1 & \text{if } A(T) < 1 \\ \int_{\{T=c\}} \varphi_* h dM / \int_{\{T=c\}} h dM & \text{if } T = c \in \{A = 1\} \\ 0 & \text{else} \end{cases}$$

with the convention that $0/0 := 0$. This does not change the power function of φ_* , because

$$\int_{\{T=c\}} \varphi_* d\mathbb{P}_\theta = C(\theta) e^{\theta c} \int_{\{T=c\}} \varphi_* h dM = C(\theta) e^{\theta c} \int_{\{T=c\}} \bar{\varphi}_* dM = \int_{\{T=c\}} \bar{\varphi}_* d\mathbb{P}_\theta$$

for any $\theta \in \Theta$ and $c \in \{A = 1\}$. Consequently, we may assume that the test φ_* has the form

$$\varphi_* = \begin{cases} 1 & \text{if } A(T) < 1 \\ 0 & \text{if } A(T) > 1 \\ \gamma(c) & \text{if } T = c \in \{A = 1\} \end{cases}$$

with numbers $\gamma_c \in [0, 1]$, $c \in \{A = 1\}$.

Suppose that $b_1 \leq 0 < b_2$ or $b_1 > 0 \geq b_2$. In this case $A(\cdot)$ would be strictly monotone, so $\mathbb{E}_\theta(\varphi_*)$ would be strictly increasing or strictly decreasing in $\theta \in \Theta$, see Theorem 7.13 (ii). But this would contradict the equation $\mathbb{E}_{\theta_1}(\varphi_*) = \mathbb{E}_{\theta_2}(\varphi_*)$. Hence $b_1, b_2 > 0$, and $A(\cdot)$ is strictly convex with $A(t) \rightarrow \infty$ as $|t| \rightarrow \infty$. Consequently, our test φ_* has the asserted form

$$\varphi_* = 1_{[T=c_1]}\gamma_1 + 1_{[T=c_2]}\gamma_2 + 1_{[c_1 < T < c_2]}$$

with real numbers $c_1 \leq c_2$ and $\gamma_1, \gamma_2 \in [0, 1]$, where $\gamma_2 = 0$ if $c_1 = c_2$.

Proof of part (ii). Let φ_* be a test as in part (i). For arbitrary fixed $\theta \in \Theta \setminus \{\theta_1, \theta_2\}$ let $d_{\theta_1} := \theta_1 - \theta < d_{\theta_2} := \theta_2 - \theta$. Then Lemma 7.25 shows that there exist constants $b_{\theta_1}, b_{\theta_2} > 0$ such that

$$A_\theta(t) := b_{\theta_1}e^{(\theta_1-\theta)t} + b_{\theta_2}e^{d(\theta_2-\theta)t}$$

satisfies

$$\begin{cases} A_\theta(c_1) = A_\theta(c_2) = 1, \\ A'_\theta(c_1) = 0 \end{cases} \quad \text{if } c_1 = c_2.$$

With $k_{\theta_j} := b_{\theta_j}C(\theta)/C(\theta_j)$, we may write

$$A_\theta(T) = k_{\theta_1} \frac{f_{\theta_1}}{f_\theta} + k_{\theta_2} \frac{f_{\theta_2}}{f_\theta} \begin{cases} > 1 & \text{if } f_\theta < k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}, \\ < 1 & \text{if } f_\theta > k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}. \end{cases}$$

Suppose first that $\theta_1 < \theta < \theta_2$. Lemma 7.25 shows that both components b_{θ_j} are strictly positive, and

$$A_\theta \begin{cases} > 1 & \text{on } \mathbb{R} \setminus [c_1, c_2], \\ < 1 & \text{on } (c_1, c_2). \end{cases}$$

Hence

$$\varphi_* = \begin{cases} 1 & \text{if } f_\theta > k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}, \\ 0 & \text{if } f_\theta < k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}. \end{cases}$$

Consequently, property (ii.1) follows from part (ii) of the generalized Neyman–Pearson lemma.

In case of $\theta < \theta_1$ or $\theta > \theta_2$, Lemma 7.25 yields the inequalities

$$A_\theta \begin{cases} < 1 & \text{on } \mathbb{R} \setminus [c_1, c_2], \\ > 1 & \text{on } (c_1, c_2). \end{cases}$$

Hence

$$1 - \varphi_* = \begin{cases} 1 & \text{if } f_\theta > k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}, \\ 0 & \text{if } f_\theta < k_{\theta_1}f_{\theta_1} + k_{\theta_2}f_{\theta_2}. \end{cases}$$

Consequently we may deduce from part (ii) of the generalized Neyman–Pearson lemma that $\mathbb{E}_\theta(\varphi) \leq \mathbb{E}_\theta(1 - \varphi_*)$ for any test φ such that $\mathbb{E}_{\theta_1}(\varphi) = \mathbb{E}_{\theta_2}(\varphi) = 1 - \alpha$. In other words, $\mathbb{E}_\theta(\varphi) \geq \mathbb{E}_\theta(\varphi_*)$ for any test φ such that $\mathbb{E}_{\theta_1}(\varphi) = \mathbb{E}_{\theta_2}(\varphi) = \alpha$. This is property (ii.2). Considering the special test $\varphi \equiv \alpha$ shows that φ_* is a level- α test of Θ_α . \square

Exercise 7.26. Let $K \subset \mathbb{R}^2$ be a convex set containing the three points $(0, 0)^\top$, $(\alpha, \alpha_1)^\top$ and $(\alpha_2, \alpha)^\top$ with real numbers $\alpha > 0$ and $\alpha_1, \alpha_2 > \alpha$. Show that $(\alpha, \alpha)^\top$ is an interior point of K .

7.5.3 Two-sided hypotheses, version 2

Version 2a. At first we consider tests of

$$\Theta_o := [\theta_1, \theta_2] \quad \text{versus} \quad \Theta_A := \Theta \setminus \Theta_o$$

with given interior points $\theta_1 < \theta_2$ of Θ . This testing problem is essentially the reverse of the testing problem in Theorem 7.24.

Theorem 7.27. (i) For any fixed $\alpha \in (0, 1)$ there exist real constants $\gamma_1, \gamma_2 \in [0, 1]$ and $c_1 \leq c_2$ (with $\gamma_2 = 0$ in case of $c_1 = c_2$) such that

$$\varphi_* := 1_{[T=c_1]}\gamma_1 + 1_{[T=c_2]}\gamma_2 + 1_{[T < c_1 \text{ or } T > c_2]}$$

is a test satisfying

$$\mathbb{E}_{\theta_1}(\varphi_*) = \mathbb{E}_{\theta_2}(\varphi_*) = \alpha.$$

(ii) A test φ_* as in part (i) has the following properties:

(ii.1)

$$\mathbb{E}_{\theta}(\varphi_*) \leq \alpha \quad \text{for all } \theta \in \Theta_o.$$

(ii.2) For any test φ satisfying $\mathbb{E}_{\theta_1}(\varphi) = \mathbb{E}_{\theta_2}(\varphi) = \alpha$,

$$\mathbb{E}_{\theta}(\varphi_*) \geq \mathbb{E}_{\theta}(\varphi) \quad \text{for all } \theta \in \Theta_A.$$

Proof of Theorem 7.27. We may apply Theorem 7.24 with $1 - \alpha$ in place of α to obtain an optimal level- $(1 - \alpha)$ test φ_{**} of $\Theta \setminus (\theta_1, \theta_2)$ versus (θ_1, θ_2) . Then the precise properties of φ_{**} provided by Theorem 7.24 imply that $\varphi_* := 1 - \varphi_{**}$ has the properties stated in Theorem 7.27. \square

Version 2b. Now we consider tests of

$$\Theta_o := \{\theta_o\} \quad \text{versus} \quad \Theta_A := \Theta \setminus \Theta_o$$

for a given interior point θ_o of Θ . Without further constraints on the tests, there exists no globally optimal level- α test of $\{\theta_o\}$. For let

$$\begin{aligned} \varphi_{\alpha,l} &= 1_{[T < k_{\alpha,l}]} + 1_{[T = k_{\alpha,l}]} \gamma_{\alpha,l}, \\ \varphi_{\alpha,r} &= 1_{[T > k_{\alpha,r}]} + 1_{[T = k_{\alpha,r}]} \gamma_{\alpha,r} \end{aligned}$$

with real constants $k_{\alpha,l}, k_{\alpha,r}$ and $\gamma_{\alpha,l}, \gamma_{\alpha,r} \in [0, 1]$ such that $\mathbb{E}_{\theta_o}(\varphi_{\alpha,l}) = \mathbb{E}_{\theta_o}(\varphi_{\alpha,r}) = \alpha$. Then by Theorem 7.13 (ii), the power functions of $\varphi_{\alpha,l}$ and $\varphi_{\alpha,r}$ are strictly decreasing and strictly increasing, respectively, and any test φ with $\mathbb{E}_{\theta_o}(\varphi) \leq \alpha$ satisfies

$$\mathbb{E}_{\theta}(\varphi) \leq \begin{cases} \mathbb{E}_{\theta}(\varphi_{\alpha,l}) & \text{if } \theta \leq \theta_o \\ \mathbb{E}_{\theta}(\varphi_{\alpha,r}) & \text{if } \theta \geq \theta_o \end{cases}$$

To obtain a unique optimal test we shall restrict our attention to tests φ such that

$$\mathbb{E}_{\theta}(\varphi) \geq \alpha = \mathbb{E}_{\theta_o}(\varphi),$$

i.e. the power of φ is nowhere lower than the power of the trivial level- α test $\varphi \equiv \alpha$.

Before going into further detail, let us mention an important property of power functions in the present setting of a one-parameter exponential family: It follows from Exercise 6.25 that for any test φ , the power function

$$\Theta \ni \theta \mapsto \mathbb{E}_\theta(\varphi)$$

is continuous on Θ and continuously differentiable on the interior of Θ with derivative

$$\begin{aligned} \frac{d}{d\theta} \mathbb{E}_\theta(\varphi) &= \frac{d}{d\theta} \frac{\int \varphi e^{\theta T} h dM}{\int e^{\theta T} h dM} \\ &= \frac{\int \varphi T e^{\theta T} h dM}{\int e^{\theta T} h dM} - \frac{\int \varphi e^{\theta T} dM \int T e^{\theta T} dM}{(\int e^{\theta T} dM)^2} \\ &= \int \varphi T d\mathbb{P}_\theta - \int \varphi d\mathbb{P}_\theta \int T d\mathbb{P}_\theta \\ &= \text{Cov}_\theta(\varphi, T). \end{aligned}$$

Hence, if φ is a test such that

$$(7.4) \quad \mathbb{E}_\theta(\varphi) \geq \alpha \geq \mathbb{E}_{\theta_o}(\varphi) \quad \text{for all } \theta \in \Theta \setminus \{\theta_o\},$$

then

$$\mathbb{E}_{\theta_o}(\varphi) = \alpha \quad \text{and} \quad \text{Cov}_{\theta_o}(\varphi, T) = 0.$$

Theorem 7.28. (i) For any fixed $\alpha \in (0, 1)$, there exist real constants $\gamma_1, \gamma_2 \in [0, 1]$ and $c_1 \leq c_2$ (with $\gamma_2 = 0$ in case of $c_1 = c_2$) such that

$$\varphi_* := 1_{[T=c_1]} \gamma_1 + 1_{[T=c_2]} \gamma_2 + 1_{[T < c_1 \text{ or } T > c_2]}$$

is a test satisfying

$$\mathbb{E}_{\theta_o}(\varphi_*) = \alpha \quad \text{and} \quad \text{Cov}_{\theta_o}(\varphi_*, T) = 0.$$

(ii) A test φ_* as in part (i) has the following property: For any test φ such that $\mathbb{E}_{\theta_o}(\varphi) = \alpha$ and $\text{Cov}_{\theta_o}(\varphi, T) = 0$,

$$\mathbb{E}_\theta(\varphi_*) \geq \mathbb{E}_\theta(\varphi) \quad \text{for all } \theta \in \Theta.$$

In particular, φ_* has property (7.4), and its power function is pointwise maximal among all tests with that property.

Proof of Theorem 7.28. Without loss of generality we may assume that $\theta_o = 0$ and $M = \mathbb{P}_0$, $h \equiv 1$, because

$$\frac{f_\theta}{f_{\theta_o}} = \frac{C(\theta)}{C(\theta_o)} \exp((\theta - \theta_o)T) \quad \text{on } \{h > 0\},$$

that means, $C(\theta)C(\theta_o)^{-1} \exp((\theta - \theta_o)T)$ is a density of \mathbb{P}_θ with respect to \mathbb{P}_{θ_o} . We may further assume without loss of generality that $\mathbb{E}_0(T) = 0$, because

$$C(\theta) \cdot \exp(\theta T) = C(\theta) \exp(\theta \mathbb{E}_0(T)) \cdot \exp(\theta(T - \mathbb{E}_0(T))).$$

Now we construct for a fixed $\theta_1 \in \Theta \setminus \{0\}$ a test φ_* maximizing $\mathbb{E}_{\theta_1}(\varphi_*)$ under the constraints that

$$\mathbb{E}_0(\varphi_*) = \int \varphi_* d\mathbb{P}_0 = \alpha \quad \text{and} \quad \mathbb{E}_0(\varphi_* T) = \int \varphi_* T d\mathbb{P}_0 = 0.$$

This may be achieved with the generalized Neyman–Pearson lemma applied to $f_1 := 1$, $f_2 := T$, $f_3 := f_{\theta_1}$ and $M := \mathbb{P}_0$. With the four tests $\varphi := 0$, $\varphi := 1$, $\varphi := 1_{[T>0]}$ and $\varphi := 1_{[T\leq 0]}$ it follows that

$$\mathcal{K}_2 := \left\{ \int \varphi \mathbf{f} d\mathbb{P}_0 : \varphi \in \mathcal{T} \right\}$$

with $\mathbf{f} := (f_1, f_2)^\top$ contains the points $(0, 0)^\top$, $(1, 0)^\top$, $(\beta, \gamma)^\top$ and $(1 - \beta, -\gamma)^\top$ with $\beta = \mathbb{P}_0(T > 0) \in (0, 1)$ and $\gamma = \int T^+ d\mathbb{P}_0 = \int T^- d\mathbb{P}_0 > 0$. Hence $(\alpha, 0)^\top$ is an interior point of \mathcal{K}_2 . Consequently there exist a test φ_* and real constants k_1, k_2 such that

$$\int \varphi_* d\mathbb{P}_0 = \alpha, \quad \int \varphi_* T d\mathbb{P}_0 = 0$$

and

$$\varphi_* = \begin{cases} 1 & \text{if } f_{\theta_1} > k_1 + k_2 T, \\ 0 & \text{if } f_{\theta_1} < k_1 + k_2 T. \end{cases}$$

Note that $f_1 = C(\theta_1)e^{\theta_1 T}$. Since $t \mapsto C(\theta_1)e^{\theta_1 t}$ is strictly convex, and since $\varphi_* \not\equiv 1$, there have to exist points $c_1 \leq c_2$ such that $C(\theta_1)e^{\theta_1 c_j} = k_1 + k_2 c_j$ for $j = 1, 2$, and

$$\varphi_* = \begin{cases} 1 & \text{if } T \in \mathbb{R} \setminus [c_1, c_2], \\ 0 & \text{if } T \in (c_1, c_2). \end{cases}$$

Finally, we may replace φ_* with

$$\bar{\varphi}_* := \begin{cases} 1 & \text{if } T \in \mathbb{R} \setminus [c_1, c_2] \\ \int_{\{T=c\}} \varphi_* d\mathbb{P}_0 / \mathbb{P}_0(T=c) & \text{if } T = c \in \{c_1, c_2\} \\ 0 & \text{else} \end{cases}$$

with the convention that $0/0 := 0$. This does not change the power function of φ_* , and then $\bar{\varphi}_*$ has the properties mentioned in part (i).

As to part (ii), consider an arbitrary fixed parameter $\theta \neq 0$. One can easily verify that there exist real constants $k_{\theta 1}, k_{\theta 2}$ such that $A_\theta(t) := C(\theta)e^{\theta t} - k_{\theta 1} - k_{\theta 2}t$ satisfies

$$\begin{cases} A_\theta(c_1) = A_\theta(c_2) = 0, \\ A'_\theta(c_1) = 0 & \text{if } c_1 = c_2. \end{cases}$$

By strict convexity of A_θ ,

$$A_\theta(t) \begin{cases} > 0 & \text{if } t \in \mathbb{R} \setminus [c_1, c_2], \\ < 0 & \text{if } t \in (c_1, c_2), \end{cases}$$

whence

$$\varphi_* = \begin{cases} 1 & \text{if } f_\theta > k_{\theta 1} + k_{\theta 2} T, \\ 0 & \text{if } f_\theta < k_{\theta 1} + k_{\theta 2} T. \end{cases}$$

Consequently, it follows from part (ii) of the generalized Neyman–Pearson lemma that $\mathbb{E}_\theta(\varphi_*) \geq \mathbb{E}_\theta(\varphi)$ for any test φ such that $\int \varphi d\mathbb{P}_0 = \alpha$ and $\int \varphi T d\mathbb{P}_0 = 0$. Taking $\varphi \equiv \alpha$ reveals that φ_* has property (7.4) and is optimal among all tests with that property. \square

Exercise 7.29. This exercise provides further details about the power function of statistical tests in a one-parameter exponential family with natural parametrization and test statistic $T : \Omega \rightarrow \mathbb{R}$.

(a) Suppose that $f : \Omega \rightarrow \mathbb{R}$ is such that

$$h_f(\theta) := \mathbb{E}_\theta(f)$$

is well-defined in \mathbb{R} for all $\theta \in \Theta$. Show that h_f is continuous on Θ and differentiable on the interior of Θ with

$$h'_f(\theta) = \mathbb{E}_\theta(fT) - \mathbb{E}_\theta(f)\mathbb{E}_\theta(T) = \text{Cov}_\theta(f, T).$$

Hint: Consider Exercise 6.25.

(b) Show that for interior points θ of Θ ,

$$h''_f(\theta) = \text{Cov}_\theta(f, T^2) - 2h'_f(\theta)\mathbb{E}_\theta(T).$$

(c) Let θ be an interior point of Θ such that $h'_f(\theta) = 0$. Show that

$$h''_f(\theta) = \mathbb{E}_\theta((T - c_1)(T - c_2)(f - \mathbb{E}_\theta(f)))$$

for arbitrary $c_1, c_2 \in \mathbb{R}$.

(d) Now consider optimal tests of two-sided hypotheses, that is

$$\varphi_* := 1_{[T=c_1]}\gamma_1 + 1_{[T=c_2]}\gamma_2 + \begin{cases} 1_{[c_1 < T < c_2]} & \text{(Type 1)} \\ 1_{[T < c_1 \text{ or } T > c_2]} & \text{(Type 2)} \end{cases}$$

with $c_1 \leq c_2$ and $\gamma_1, \gamma_2 \in [0, 1]$. Show that

$$(\varphi_* - \mathbb{E}_\theta(\varphi_*))(T - c_1)(T - c_2) \begin{cases} \leq 0 & \text{(Type 1)} \\ \geq 0 & \text{(Type 2)} \end{cases}$$

with strict inequality in case of $T \notin \{c_1, c_2\}$. Deduce from this and part (c), that in case of $h'_{\varphi_*}(\theta) = 0$,

$$h''_{\varphi_*}(\theta) \begin{cases} < 0 & \text{(type 1)} \\ > 0 & \text{(type 2)} \end{cases}$$

unless \mathbb{P}_θ^T is concentrated on $\{c_1, c_2\}$.

7.5.4 Summary and some first applications

For notational convenience, we formulated and derived the previous results for one-parameter exponential families with “natural parametrization”. That means, the exponential term of the density f_θ contains the product of the test statistic T with the parameter θ . In the examples we looked at in Section 7.5.1, this necessitated a transformation of the original parameters. To get a more complete picture, let us summarize the main results of this section in terms of the original parametrization.

Definition 7.30 (One-parameter exponential family). A one-parameter exponential family is a statistical experiment

$$(\Omega, \mathcal{A}, (\mathbb{P}_\lambda)_{\lambda \in \Lambda})$$

of the following form: The parameter space Λ is a real interval. For a σ -finite measure M on (Ω, \mathcal{A}) and measurable functions $h : \Omega \rightarrow [0, \infty)$, $T : \Omega \rightarrow \mathbb{R}$, each distribution \mathbb{P}_λ has density $f_\lambda = d\mathbb{P}_\lambda/dM$ given by

$$f_\lambda(\omega) = C(\lambda)h(\omega)\exp(\theta(\lambda)T(\omega)),$$

where $\theta : \Lambda \rightarrow \mathbb{R}$ is a differentiable mapping with $\theta' > 0$ on Λ or $\theta' < 0$ on Λ . Moreover, $M(h > 0 \text{ and } T \neq c) > 0$ for any real constant, so $\mathbb{P}_{\lambda_1} \neq \mathbb{P}_{\lambda_2}$ whenever $\lambda_1 \neq \lambda_2$.

The parameter $\theta(\lambda)$ in the previous definition is called “natural parameter”. The set $\Theta = \theta(\Lambda)$ is a subset of the “natural parameter space”

$$\Theta_{\text{nat}} := \left\{ \theta \in \mathbb{R} : \int h e^{\theta T} dM < \infty \right\}.$$

Now we consider tests φ_* of one of the following types:

$$(7.5) \quad \varphi_* = 1_{[T=c_1]}\gamma_1 + 1_{[c_1 < T < c_2]} + 1_{[T=c_2]}\gamma_2,$$

$$(7.6) \quad \varphi_* = 1_{[T=c_1]}\gamma_1 + 1_{[T=c_2]}\gamma_2 + 1_{[T < c_1 \text{ or } T > c_2]},$$

where $c_1 \leq c_2$ and $\gamma_1, \gamma_2 \in [0, 1]$ (with $\gamma_2 = 0$ if $c_1 = c_2$).

If M^T is continuous in the sense that $M(T = c) = 0$ for any $c \in \mathbb{R}$, it suffices to consider tests φ_* of the following type:

$$(7.7) \quad \varphi_* = 1_{[c_1 \leq T \leq c_2]},$$

$$(7.8) \quad \varphi_* = 1_{[T \leq c_1 \text{ or } T \geq c_2]},$$

where $c_1 < c_2$.

Two-sided test, version 1. For given points $\lambda_1 < \lambda_2$ in Λ , an optimal level- α test of

$$\Lambda \setminus (\lambda_1, \lambda_2) \quad \text{versus} \quad (\lambda_1, \lambda_2)$$

is given by (7.5) or (7.7), provided that

$$\mathbb{E}_{\lambda_1}(\varphi_*) = \alpha = \mathbb{E}_{\lambda_2}(\varphi_*).$$

Two-sided test, version 2a. For given interior points $\lambda_1 < \lambda_2$ of Λ , an optimal level- α test of

$$[\lambda_1, \lambda_2] \quad \text{versus} \quad \Lambda \setminus [\lambda_1, \lambda_2]$$

with exact power α at λ_1 and λ_2 is given by (7.6) or (7.8), provided that

$$\mathbb{E}_{\lambda_1}(\varphi_*) = \alpha = \mathbb{E}_{\lambda_2}(\varphi_*).$$

Two-sided test, version 2b. For a given interior point λ_o of Λ , an optimal level- α test of

$$\{\lambda_o\} \quad \text{versus} \quad \Lambda \setminus \{\lambda_o\}$$

with power function bounded from below by α is given by (7.6) or (7.8), provided that

$$\mathbb{E}_{\lambda_o}(\varphi_*) = \alpha \quad \text{and} \quad \left. \frac{d}{d\lambda} \right|_{\lambda=\lambda_o} \mathbb{E}_{\lambda}(\varphi_*) = 0.$$

Example 7.11 (Gaussian location family, cont.) Suppose we observe a random vector $\mathbf{X} \in \mathbb{R}^n$ with distribution $\mathbb{P}_{\mu} := \mathcal{N}(\mu, \sigma^2)^{\otimes n}$ for a given $\sigma > 0$ and some unknown $\mu \in \mathbb{R}$. As shown before, the model $(\mathcal{N}(\mu, \sigma^2)^{\otimes n})_{\mu \in \mathbb{R}}$ is a one-parameter exponential family with natural parameter $\theta(\mu) = n\mu/\sigma^2$ and test statistic $T(\mathbf{X}) := \bar{X}$. For any fixed μ_o , an optimal level- α test of

$$\{\mu_o\} \quad \text{versus} \quad \mathbb{R} \setminus \{\mu_o\}$$

is given by

$$\varphi_*(\mathbf{X}) := \begin{cases} 1 & \text{if } |\bar{X} - \mu_o| \geq \Phi^{-1}(1 - \alpha/2)\tau, \\ 0 & \text{else,} \end{cases}$$

where $\tau := \sigma/\sqrt{n}$ is the standard deviation of \bar{X} . This follows from the fact that φ_* is of type (7.8), and with $Z := (\bar{X} - \mu)/\tau \sim \mathcal{N}(0, 1)$, the power function of φ_* is given by

$$\begin{aligned} \mathbb{E}_{\mu}(\varphi_*) &= \mathbb{P}\left(Z \geq \frac{\mu_o - \mu}{\tau} + \Phi^{-1}(1 - \alpha/2)\right) + \mathbb{P}\left(Z \leq \frac{\mu_o - \mu}{\tau} - \Phi^{-1}(1 - \alpha/2)\right) \\ &= 1 - \Phi\left(\frac{\mu_o - \mu}{\tau} + \Phi^{-1}(1 - \alpha/2)\right) + \Phi\left(\frac{\mu_o - \mu}{\sigma/\sqrt{n}} - \Phi^{-1}(1 - \alpha/2)\right) \\ &= \Phi\left(\Phi^{-1}(\alpha/2) + \frac{\mu - \mu_o}{\tau}\right) + \Phi\left(\Phi^{-1}(\alpha/2) - \frac{\mu - \mu_o}{\tau}\right). \end{aligned}$$

This equals α for $\mu = \mu_o$, and it is an even function of $\mu - \mu_o$, so the derivative of the power function at μ_o equals 0.

Suppose we want to verify the working hypothesis that

$$|\mu - \mu_o| < \delta$$

for given numbers $\mu_o \in \mathbb{R}$ and $\delta > 0$. This corresponds to Version 1 of a two-sided test with boundary parameters $\lambda_1 = \mu_o - \delta$ and $\lambda_2 = \mu_o + \delta$. By symmetry reasons, a possible ansatz for φ_* would be

$$\varphi_*(\mathbf{X}) := \begin{cases} 1 & \text{if } |\bar{X} - \mu_o| \leq c_{\alpha}\tau \\ 0 & \text{else} \end{cases}$$

for a suitable constant $c_{\alpha} > 0$. Indeed, this defines a test of type (7.7), and with Z as above we may write

$$\begin{aligned} \mathbb{E}_{\mu_o \pm \delta}(\varphi_*) &= \mathbb{P}_{\pm\delta}(|\bar{X}| \leq c_{\alpha}\tau) = \mathbb{P}(|\pm\delta/\tau + Z| \leq c_{\alpha}) \\ &= \Phi(c_{\alpha} \mp \delta/\tau) - \Phi(-c_{\alpha} \mp \delta/\tau) \\ &= \Phi(c_{\alpha} + \delta/\tau) + \Phi(c_{\alpha} - \delta/\tau) - 1. \end{aligned}$$

This is obviously a continuous and strictly increasing function of c_α with value 0 for $c_\alpha = 0$ and limit 1 as $c_\alpha \rightarrow \infty$. Hence there exists a unique $c_\alpha(\delta/\tau) > 0$ such that

$$\mathbb{E}_{\mu_o \pm \delta}(\varphi_*) = \alpha.$$

The precise value of $c_\alpha(\delta/\tau)$ has to be computed numerically.

If we let $\delta \downarrow 0$, we obtain $c_\alpha(0) = \Phi^{-1}((1 + \alpha)/2)$. Indeed, an optimal level- α test of

$$\mathbb{R} \setminus \{\mu_o\} \quad \text{versus} \quad \{\mu_o\}$$

rejects the null hypothesis if

$$|\bar{X} - \mu_o| \leq \tau \Phi^{-1}((1 + \alpha)/2).$$

In this case, we may claim with confidence $1 - \alpha$ that $\mu = \mu_o$. This sounds almost miraculous, but note that the inequality $|\bar{X} - \mu_o| \leq \tau \Phi^{-1}((1 + \alpha)/2)$ occurs with probability at most α . Instead of trying to prove that $\mu = \mu_o$, one should rather compute an upper $(1 - \alpha)$ -confidence bound for $|\mu - \mu_o|$, see later.

Exercise 7.31. Let $X \sim \text{Gamma}(a, b)$ with shape parameter $a > 0$ and scale parameter $b > 0$, i.e. X has density $f_{a,b}$ with respect to Lebesgue measure on $(0, \infty)$, where

$$f_{a,b}(x) = \Gamma(a)^{-1} b^{-1} (x/b)^{a-1} e^{-x/b}.$$

(a) Suppose that $a > 0$ is given but $b > 0$ is unknown. Verify that $(\text{Gamma}(a, b))_{b>0}$ corresponds to a one-parameter exponential family with test statistic $T(X) = X$.

(b) Suppose that $b > 0$ is given but $a > 0$ is unknown. Verify that $(\text{Gamma}(a, b))_{a>0}$ corresponds to a one-parameter exponential family with test statistic $T(X) = \log(X)$.

(c) Assuming that $a > 0$ is given but $b > 0$ is an unknown parameter in $(0, \infty)$, determine an optimal level- α test of

$$\{1\} \quad \text{versus} \quad (0, \infty) \setminus \{1\}.$$

(d) Modify your test in part (b) to become an optimal test of

$$\{b_o\} \quad \text{versus} \quad (0, \infty) \setminus \{b_o\}.$$

for arbitrary fixed $b_o > 0$.

Exercise 7.32. Let $X \sim \text{Bin}(n, p)$ with given $n \in \mathbb{N}$ and unknown $p \in [0, 1]$.

(a) Fix small numbers $\alpha \in (0, 1)$ and $\delta \in (0, 0.5)$. Construct a statistical procedure to verify with given confidence $1 - \alpha$ that $|p - 0.5| < \delta$.

(b) How large should n be such that for a given small $\alpha' \in (0, 1)$, this conclusion is drawn with probability at least $1 - \alpha'$ in case of $p = 0.5$? Give a numerical answer to this question in case of $\alpha = 0.05$, $\alpha' = 0.3$ and $\delta = 0.1$.

7.6 Tests and Confidence regions

For any statistical model $(\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ there is a close relationship between tests and confidence regions. Let us first clarify what we mean by confidence regions.

Definition 7.33 (Randomized confidence region). A *(randomized) confidence region* is a mapping $C : \Omega \times \Theta \rightarrow [0, 1]$ such that for any fixed $\theta \in \Theta$, the mapping $C(\cdot, \theta)$ is \mathcal{A} -measurable.

Suppose that for a given test level $\alpha \in (0, 1)$,

$$\int C(\omega, \theta) \mathbb{P}_\theta(d\omega) \geq 1 - \alpha \quad \text{for arbitrary } \theta \in \Theta.$$

Then C is called a *(randomized) confidence region with confidence level $1 - \alpha$* , or shortly: a *$(1 - \alpha)$ -confidence region*.

Interpretation, and the meaning of confidence. Suppose that we observe a random data set $X \in \Omega$ with distribution \mathbb{P}_θ , where $\theta \in \Theta$ is unknown. Let U be a random variable with uniform distribution on $[0, 1]$ and independent from X . Then we claim that θ is contained in the set

$$\mathcal{C}(X, U) := \{\theta \in \Theta : C(X, \theta) \geq U\} \subset \Theta.$$

In case of C being non-randomized, i.e. C taking only values in $\{0, 1\}$, we don't need the extra random variable U , because the set $\mathcal{C}(X, U)$ equals

$$\mathcal{C}(X) := \{\theta \in \Theta : C(X, \theta) = 1\} \subset \Theta$$

almost surely.

If C is a $(1 - \alpha)$ -confidence region, then

$$\left. \begin{array}{l} \mathbb{P}_\theta(\theta \in \mathcal{C}(X)) \\ \mathbb{P}_\theta(\theta \in \mathcal{C}(X, U)) \end{array} \right\} \geq 1 - \alpha \quad \text{for all } \theta \in \Theta.$$

That means, the confidence region covers the unknown true parameter with probability at least $1 - \alpha$. This statement involves the *random variable* X (and U , if needed) and is true *prior* to observing X (and generating U).

Once we have observed X (and generated U , if needed), the claim that $\mathcal{C}(X)$ or $\mathcal{C}(X, U)$ contains the unknown true parameter is simply true or false. Hence it would be ridiculous to say that “with probability $1 - \alpha$, the confidence region $\mathcal{C}(X)$ or $\mathcal{C}(X, U)$ contains θ ”. Instead, one may claim with *confidence* $1 - \alpha$ that $\theta \in \mathcal{C}(X)$ or $\theta \in \mathcal{C}(X, U)$. This formulation indicates that for a specific observed data set, the claim is simply true or false, but we use a procedure which leads to a correct statement in at least $(1 - \alpha) \cdot 100$ percent of applications in the long run.

Duality between confidence regions and tests of one-point hypotheses. If $C : \Omega \times \Theta \rightarrow [0, 1]$ is a (randomized) confidence region, then for any $\theta_o \in \Theta$, a test of $\{\theta_o\}$ is given by $\varphi(\cdot, \theta_o) := 1 - C(\cdot, \theta_o)$. On the other hand, if for any $\theta_o \in \Theta$ we have defined a test $\varphi(\cdot, \theta_o)$ of $\{\theta_o\}$, then $C(\omega, \theta) := 1 - \varphi(\omega, \theta)$ defines a confidence region C . The confidence region C has confidence level $1 - \alpha$ if and only if each test $\varphi(\cdot, \theta_o)$ is a level- α test of $\{\theta_o\}$.

Example 7.34 (Simple Gaussian shift model). Suppose we observe $X \sim \mathbb{P}_\mu = \mathcal{N}(\mu, 1)$ for some unknown parameter $\mu \in \mathbb{R}$. One can easily verify that for any hypothetical parameter $\mu_o \in \mathbb{R}$, an optimal level- α test of $\{\mu_o\}$ with power function at least α everywhere is given by

$$\varphi(X, \mu_o) := 1_{[|X - \mu_o| > \Phi^{-1}(1 - \alpha/2)]}.$$

This leads to the $(1 - \alpha)$ -confidence region C given by

$$C(X, \mu) = 1_{[|X - \mu| \leq \Phi^{-1}(1 - \alpha/2)]} = 1_{[\mu \in [X \pm \Phi^{-1}(1 - \alpha/2)]]},$$

i.e. $\mathcal{C}(X) = [X \pm \Phi^{-1}(1 - \alpha/2)]$.

Exercise 7.35. As in Exercise 7.31, suppose we observe $X \sim \text{Gamma}(a, b)$ with given shape parameter $a > 0$ and unknown scale parameter $b > 0$. An ansatz for a confidence interval for b is

$$\mathcal{C}(X) := [X/\kappa_2, X/\kappa_1]$$

with constants $0 < \kappa_1 < \kappa_2$.

(a) Show that for any choice of (κ_1, κ_2) , the coverage probability $\mathbb{P}_p(p \in \mathcal{C}(X))$ is constant in $b > 0$. Then characterize the set of all pairs (κ_1, κ_2) such that the confidence level is exactly equal to $1 - \alpha$.

(b) A potential measure for the size of $\mathcal{C}(X)$ is the ratio of its upper and lower boundary,

$$\frac{X/\kappa_1}{X/\kappa_2} = \frac{\kappa_2}{\kappa_1},$$

or the logarithm thereof. Determine the unique pair (κ_1, κ_2) minimizing this size. (The solution is characterized by some equation which could be solved numerically for specific α .) Show also that this pair satisfies $\kappa_1 < a < \kappa_2$.

(c) What is the relation of the solution in part (b) to the optimal test in Exercise 7.31 (b)?

Duality between confidence regions and tests of composite hypotheses. Sometimes we are not interested in a confidence region for the unknown true parameter θ but only in an upper bound for $g(\theta)$ with a given function

$$g : \Theta \rightarrow \mathbb{R}.$$

For instance, if (Θ, d) is a metric space, we are sometimes interested in the distance between the unknown true parameter and some given point $\theta_o \in \Theta$, so $g(\theta) := d(\theta, \theta_o)$.

An upper confidence bound $b_\alpha(X)$ for $g(\theta)$ corresponds to a measurable function $b_\alpha : \Omega \rightarrow \mathbb{R}$ such that

$$\mathbb{P}_\theta(b_\alpha \geq g(\theta)) \geq 1 - \alpha \quad \text{for all } \theta \in \Theta.$$

Such a confidence bound gives rise to the level- α test

$$\varphi_\alpha(X, \delta) := 1_{[b_\alpha(X) < \delta]}$$

of the null hypothesis

$$\Theta(\delta) := \{\theta \in \Theta : g(\theta) \geq \delta\}$$

for arbitrary numbers $\delta \geq 0$. Indeed, if $g(\theta) \geq \delta$, then

$$\mathbb{P}_\theta(\varphi_\alpha(\cdot, \delta) = 1) = \mathbb{P}_\theta(b_\alpha < \delta) \leq \mathbb{P}_\theta(b_\alpha < g(\theta)) \leq \alpha.$$

On the other hand, suppose that for any number $\delta \geq 0$, we have constructed a level- α test $\varphi_\alpha(\cdot, \delta) : \Omega \rightarrow \{0, 1\}$ of the null hypothesis $\Theta(\delta)$. Then

$$b_\alpha(X) := \sup(\{0\} \cup \{\delta \geq 0 : \varphi_\alpha(X, \delta) = 0\})$$

defines an upper $(1 - \alpha)$ -confidence bound for $g(\theta)$, because for any $\theta \in \Theta$ and $\delta := g(\theta)$,

$$\mathbb{P}_\theta(g(\theta) \leq b_\alpha) \geq \mathbb{P}_\theta(\varphi_\alpha(\cdot, \delta) = 0) = 1 - \mathbb{P}_\theta(\varphi_\alpha(\cdot, \delta) = 1) \geq 1 - \alpha.$$

Finding an “optimal” upper confidence bound b_α may be interpreted as finding optimal level- α tests $\varphi_\alpha(\cdot, \delta)$, $\delta \geq 0$, of the null hypotheses $\Theta(\delta)$. Ideally, the function $\varphi_\alpha(\cdot, \cdot)$ is even non-decreasing in its second argument.

Example 7.34 (Simple Gaussian shift model, cont.). As shown in the next exercise, for any given $\mu_o \in \mathbb{R}$, a simple upper $(1 - \alpha)$ -confidence bound for $|\mu - \mu_o|$ is given by $|X - \mu_o| + \Phi^{-1}(1 - \alpha)$. But if we think about the duality of tests and confidence regions, a good upper $(1 - \alpha)$ -confidence bound $b_\alpha(X)$ for $|\mu - \mu_o|$ should satisfy the following condition: For any $\delta \geq 0$,

$$\varphi_\alpha(X, \delta) := 1_{[b_\alpha(X) < \delta]}$$

defines an optimal level- α test of

$$\mathbb{R} \setminus [\mu_o \pm \delta] \quad \text{versus} \quad [\mu_o \pm \delta].$$

This is essentially Version 1 of our two-sided testing problem, except that the alternative hypothesis is chosen to be a closed rather than an open interval. We know already a solution for this testing problem: An optimal level- α test is given by

$$\varphi_\alpha(X, \delta) := 1_{[|X - \mu_o| < c_\alpha(\delta)]}$$

where $c_\alpha(\delta) > 0$ is the unique number such that

$$\Phi(c_\alpha(\delta) + \delta) + \Phi(c_\alpha(\delta) - \delta) = 1 + \alpha.$$

Recall that $c_\alpha(0) = \Phi^{-1}((1 + \alpha)/2)$, but for $\delta > 0$, there is no simple formula for $c(\delta)$. An optimal upper confidence bound for $|\mu - \mu_o|$ is given by

$$\begin{aligned} b_\alpha(X) &:= \max(\{0\} \cup \{\delta \geq 0 : \varphi_\alpha(X, \delta) = 0\}) \\ &= \max(\{0\} \cup \{\delta \geq 0 : c(\delta) \leq |X - \mu_o|\}) \\ &= \max(\{0\} \cup \{\delta \geq 0 : \Phi(|X - \mu_o| + \delta) + \Phi(|X - \mu_o| - \delta) \geq 1 + \alpha\}). \end{aligned}$$

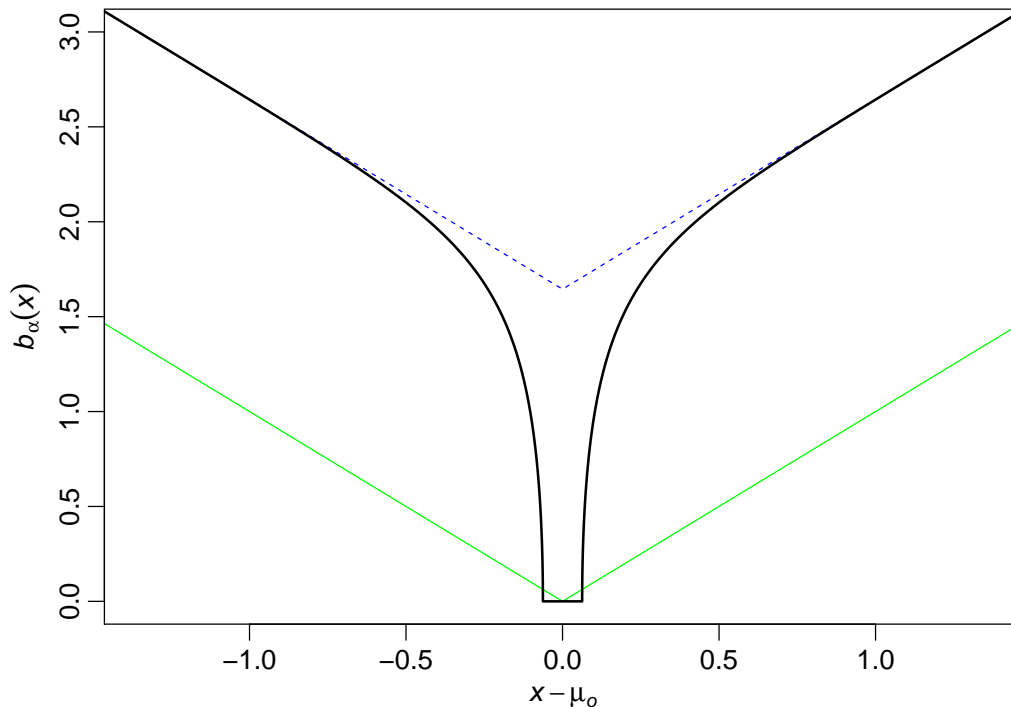


Figure 7.1: Optimal (black) and ad hoc (blue) upper 95%-confidence bound for $|\mu|$ in the simple Gaussian shift model.

If $|X - \mu_o| \leq \Phi^{-1}((1 + \alpha)/2)$, then $b_\alpha(X) = 0$. Otherwise, $b_\alpha(X)$ is the unique solution $\delta > 0$ of the equation

$$\Phi(|X - \mu_o| + \delta) + \Phi(|X - \mu_o| - \delta) = 1 + \alpha,$$

which can be computed numerically.

Figure 7.1 depicts the optimal upper 95%-confidence bound (black curve) for $|\mu - \mu_o|$ and the simple upper bound $|X - \mu_o| + \Phi^{-1}(1 - \alpha)$ (blue broken line).

Exercise 7.36 (Simple Gaussian shift model). Let $X \sim \mathbb{P}_\mu := \mathcal{N}(\mu, 1)$ for some unknown parameter $\mu \in \mathbb{R}$. A standard $(1 - \alpha)$ -confidence interval for μ is given by $[X \pm \Phi^{-1}(1 - \alpha/2)]$.

(a) Show that for any given $\mu_o \in \mathbb{R}$,

$$\mathcal{C}_\alpha(X) := [\min(X - \Phi^{-1}(1 - \alpha), \mu_o), \max(X + \Phi^{-1}(1 - \alpha), \mu_o)]$$

is also a $(1 - \alpha)$ -confidence interval for μ , i.e.

$$\mathbb{P}_\mu(\mu \in \mathcal{C}_\alpha) \geq 1 - \alpha \quad \text{for all } \mu \in \mathbb{R}.$$

(b) Show that

$$b_\alpha(X) := |X - \mu_o| + \Phi^{-1}(1 - \alpha)$$

is an upper $(1 - \alpha)$ -confidence bound for $|\mu - \mu_o|$, i.e.

$$\mathbb{P}_\mu(|\mu - \mu_o| \leq b_\alpha) \geq 1 - \alpha \quad \text{for all } \mu \in \mathbb{R}.$$

Exercise 7.37. Suppose we observe $X \sim \text{Bin}(n, p)$ with given $n \in \mathbb{N}$ and unknown $p \in [0, 1]$.

(a) Show that for arbitrary integers $0 \leq c_1 \leq c_2 \leq n$ with $c_2 - c_1 < n$,

$$[0, 1] \ni p \mapsto \log \mathbb{P}_p(c_1 \leq X \leq c_2) \in [-\infty, 0]$$

is strictly concave.

(b) Describe an optimal upper $(1 - \alpha)$ -confidence bound $b_\alpha(X)$ for $|p - 0.5|$.

Hint: Consider $\mathbb{P}_p(k \leq X \leq n - k)$ for $k = 0, 1, \dots, \lfloor n/2 \rfloor$.

(c) Write a computer program to compute the confidence bound $b_\alpha(X)$ in part (b).

Chapter 8

Decision Problems and Procedures, Sufficiency and Completeness

In the present chapter we introduce some fundamental concepts and results from statistical decision theory. Estimation and testing problems may be viewed as special cases of *decision problems*, while point estimators and statistical tests are corresponding *decision procedures*.

8.1 Decision Problems and Procedures

As in the previous chapters, we consider a statistical experiment

$$\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta}).$$

If there is no doubt about the sample space (Ω, \mathcal{A}) , we just write $\mathcal{E} = (\mathbb{P}_\theta)_{\theta \in \Theta}$.

Decision spaces. A *decision space* is a measurable space $(\mathbb{V}, \mathcal{B})$ representing the possible conclusions we could draw about the unknown true parameter θ .

Loss functions. A *loss function* is a mapping

$$L : \mathbb{V} \times \Theta \rightarrow (-\infty, \infty]$$

such that $L(\cdot, \theta)$ is \mathcal{B} -measurable for any fixed $\theta \in \Theta$. Here $L(v, \theta)$ quantifies the loss (e.g. the costs) when drawing the conclusion $v \in \mathbb{V}$ while the true parameter equals $\theta \in \Theta$.

Example 8.1 (Point estimation). Consider a given mapping $g : \Theta \rightarrow \mathbb{R}^q$. For instance, for the statistical experiment $\mathcal{E} = (\mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma) \in \mathbb{R} \times (0, \infty)}$, one could think about $g(\mu, \sigma) := \mu$ or $g(\mu, \sigma) := \sigma$. Our decision could be a guess $v \in \mathbb{R}^q$ for the unknown true value $g(\theta)$. A potential loss function would be given by

$$L(v, \theta) := \|v - g(\theta)\|^r$$

for some norm $\|\cdot\|$ on \mathbb{R}^q and some exponent $r > 0$. In case of $r \geq 1$, this loss function is a special case of the more general loss function L given by

$$L(v, \theta) := \psi(v - g(\theta))$$

for some convex function $\psi : \mathbb{R}^q \rightarrow [0, \infty)$ such that $\psi(v) \rightarrow \infty$ as $\|v\| \rightarrow \infty$.

Example 8.2 (Statistical tests). Suppose we have split Θ into two disjoint sets Θ_0 and Θ_1 . The question is whether the unknown true parameter θ belongs to Θ_0 or to Θ_1 . Hence the two potential decisions would be 0 and 1, representing the claims that $\theta \in \Theta_0$ or $\theta \in \Theta_1$, respectively. A potential loss function would be the indicator of a wrong conclusion, i.e.

$$L(v, \theta) := 1_{[v=1, \theta \in \Theta_0]} + 1_{[v=0, \theta \in \Theta_1]}.$$

More generally, one could specify costs $\lambda_j > 0$ for an error of the j -th kind and set

$$(8.1) \quad L(v, \theta) := 1_{[v=1, \theta \in \Theta_0]} \lambda_1 + 1_{[v=0, \theta \in \Theta_1]} \lambda_2.$$

Decision problems. A triplet $(\mathcal{E}, (\mathbb{V}, \mathcal{B}), L)$, consisting of a statistical experiment, a decision space and a loss function is called a *decision problem*. If there is a standard σ -field \mathcal{B} on \mathbb{V} , we just write $(\mathcal{E}, \mathbb{V}, L)$.

Decision procedures. A *non-randomized decision procedure* is a measurable mapping $\rho : \Omega \rightarrow \mathbb{V}$. That means, if we observe $X \sim \mathbb{P}_\theta$ with unknown $\theta \in \Theta$, then we draw the conclusion $\rho(X) \in \mathbb{V}$ about θ .

More generally, a *decision procedure* is a stochastic kernel ρ from (Ω, \mathcal{A}) to $(\mathbb{V}, \mathcal{B})$. That means,

$$\rho : \Omega \times \mathcal{B} \rightarrow [0, 1]$$

is a mapping such that

$$\begin{aligned} &\text{for any } \omega \in \Omega, \quad \rho(\omega, \cdot) \text{ is a probability measure on } (\mathbb{V}, \mathcal{B}), \\ &\text{for any } B \in \mathcal{B}, \quad \rho(\cdot, B) \text{ is } \mathcal{A}\text{-measurable on } \Omega. \end{aligned}$$

Now the interpretation is that having observed $X \sim \mathbb{P}_\theta$, we draw a random conclusion about θ from the probability measure $\rho(X, \cdot)$.

With slightly ambiguous notation, a non-randomized decision procedure $\rho : \Omega \rightarrow \mathbb{V}$ corresponds to the stochastic kernel $\rho(\cdot, \cdot)$ with

$$\rho(\omega, \cdot) := \delta_{\rho(\omega)}.$$

Risk functions. The performance of a decision procedure ρ is quantified by its risk function $R(\rho, \cdot) : \Theta \rightarrow [-\infty, \infty]$,

$$R(\rho, \theta) := \int_{\Omega} \int_{\mathbb{V}} L(v, \theta) \rho(\omega, dv) \mathbb{P}_\theta(d\omega).$$

In our explicit examples, the loss functions L are non-negative, so $R(\rho, \theta)$ is well-defined in $[0, \infty]$.

Example 8.1 (Point estimation, continued). A point estimator for $g(\theta)$ is a measurable mapping $\hat{g} : \Omega \rightarrow \mathbb{R}^q$. This corresponds to the (non-randomized) decision procedure

$$\rho(\omega, \cdot) := \delta_{\hat{g}(\omega)}.$$

Its risk function is given by

$$R(\hat{g}, \theta) = \int_{\Omega} \|\hat{g} - g(\theta)\|^r d\mathbb{P}_{\theta} = \mathbb{E}_{\theta}(\|\hat{g} - g(\theta)\|^r)$$

or

$$R(\hat{g}, \theta) = \int_{\Omega} \psi(\hat{g} - g(\theta)) d\mathbb{P}_{\theta} = \mathbb{E}_{\theta} \psi(\|\hat{g} - g(\theta)\).$$

Exercise 8.3 (De-randomisation for point estimation). This exercise shows that in the context of point estimation, it is often sufficient to consider non-randomized decision procedures, i.e. simple point estimators. Let $L(v, \theta) = \psi(v - g(\theta))$ with a convex function $\psi : \mathbb{R}^q \rightarrow \mathbb{R}$ such that $\psi(v) \rightarrow \infty$ as $\|v\| \rightarrow \infty$. Show that for any decision procedure ρ , i.e. a stochastic kernel from (Ω, \mathcal{A}) to $(\mathbb{R}^d, \text{Borel}(\mathbb{R}^d))$, there exists a point estimator $\hat{g} : \Omega \rightarrow \mathbb{R}^d$ such that

$$R(\hat{g}, \theta) = \int_{\Omega} L(\hat{g}(\omega), \theta) \mathbb{P}_{\theta}(d\omega)$$

is no larger than

$$R(\rho, \theta) = \int_{\Omega} \int_{\mathbb{R}^d} L(v, \theta) \rho(\omega, dv) \mathbb{P}_{\theta}(d\omega).$$

Example 8.2 (Statistical tests, continued). With the decision space $\mathbb{V} = \{0, 1\}$, any decision procedure ρ may be written as

$$\rho(\omega, \cdot) = (1 - \varphi(\omega))\delta_0 + \varphi(\omega)\delta_1$$

for some measurable function $\varphi : \Omega \rightarrow [0, 1]$, i.e. a test on Ω . With the power function $\theta \mapsto \mathbb{E}_{\theta}(\varphi)$ of φ and the loss function L in (8.1),

$$R(\varphi, \theta) = 1_{[\theta \in \Theta_0]} \lambda_1 \mathbb{E}_{\theta}(\varphi) + 1_{[\theta \in \Theta_1]} \lambda_2 (1 - \mathbb{E}_{\theta}(\varphi)).$$

Bayes risks. Suppose that Θ itself is equipped with a σ -field \mathcal{C} , and suppose that the loss function $L : \Theta \times \mathbb{V} \rightarrow (-\infty, \infty]$ is $\mathcal{C} \otimes \mathcal{B}$ -measurable. Further suppose that $\theta \mapsto \mathbb{P}_{\theta}(A)$ is \mathcal{C} -measurable for any fixed $A \in \mathcal{A}$. Then we may view $(\mathbb{P}_{\theta})_{\theta \in \Theta}$ as a stochastic kernel, too, and consider the following Bayesian model: Let Π be a probability distribution on (Θ, \mathcal{C}) , a so-called *prior (distribution)*. One could imagine “mother nature” choosing a parameter $\theta \in \Theta$ from that distribution Π , and then, conditional on θ , we observe a random variable $X \sim \mathbb{P}_{\theta}$. The Bayes risk of a decision procedure ρ for the prior Π is defined as

$$R(\rho, \Pi) := \int_{\Theta} R(\rho, \theta) \Pi(d\theta) = \int_{\Theta} \int_{\Omega} \int_{\mathbb{V}} L(\theta, v) \rho(\omega, dv) \mathbb{P}_{\theta}(d\omega) \Pi(d\theta).$$

8.2 Some Optimality Concepts and Results

Let $(\mathcal{E}, (\mathbb{V}, \mathcal{B}), L)$ be a given decision problem. Our goal is to devise decision procedures ρ with low risks $R(\rho, \theta)$. Typically there is no “free lunch”: If ρ has very small risk $R(\rho, \theta_1)$ for some parameter $\theta_1 \in \Theta$, it will often have rather large risk $R(\rho, \theta_2)$ for some other parameter $\theta_2 \in \Theta$.

Minimax-optimality. A decision procedure ρ_* is called *minimax-optimal*, if

$$\sup_{\theta \in \Theta} R(\rho_*, \theta) = \min_{\rho} \sup_{\theta \in \Theta} R(\rho, \theta).$$

Throughout this chapter, “ \min_{ρ} ” and “ \inf_{ρ} ” stand for the minimum and infimum, respectively, over all decision procedures ρ .

Admissibility. A decision procedure ρ_* is called *admissible*, if there exists no decision procedure ρ satisfying

$$R(\rho, \theta) \leq R(\rho_*, \theta) \quad \text{for all } \theta \in \Theta$$

and

$$R(\rho, \theta_o) < R(\rho_*, \theta_o) \quad \text{for at least one } \theta_o \in \Theta.$$

Bayes-optimality. For a given prior Π on Θ , a decision procedure ρ_* is called *Bayes-optimal* for this prior Π , if

$$R(\rho_*, \Pi) = \min_{\rho} R(\rho, \Pi).$$

Least favourable priors. A prior Π_* is called least favourable, if

$$\inf_{\rho} R(\rho, \Pi_*) = \max_{\Pi} \inf_{\rho} R(\rho, \Pi),$$

where “ \max_{Π} ” stands for the maximum over all priors Π on Θ .

Here are three simple results establishing minimaxity, Bayes-optimality and admissibility of decision procedures.

Lemma 8.4. *Let Π_* be a prior on Θ , and let ρ_* be a Bayes-optimal decision procedure for Π_* . Suppose further that*

$$R(\rho_*, \Pi_*) = \sup_{\theta \in \Theta} R(\rho_*, \theta).$$

Then ρ_ is minimax-optimal, and Π_* is a least favourable prior.*

Proof of Lemma 8.4. For any decision procedure ρ ,

$$\sup_{\theta \in \Theta} R(\rho, \theta) \geq R(\rho, \Pi_*),$$

and by our assumptions on ρ_* ,

$$R(\rho, \Pi_*) \geq R(\rho_*, \Pi_*) = \sup_{\theta \in \Theta} R(\rho_*, \theta).$$

Hence ρ_* is minimax-optimal.

For any prior Π on Θ ,

$$\inf_{\rho} R(\rho, \Pi) \leq R(\rho_*, \Pi) \leq \sup_{\theta \in \Theta} R(\rho_*, \theta) = R(\rho_*, \Pi_*) = \inf_{\rho} R(\rho, \Pi_*)$$

by assumption. Hence Π_* is least favourable. \square

Lemma 8.5. *Let ρ_* be a decision procedure such that for a sequence $(\Pi_k)_{k \geq 1}$ of priors,*

$$\sup_{\theta \in \Theta} R(\rho_*, \theta) = \lim_{k \rightarrow \infty} \inf_{\rho} R(\rho, \Pi_k).$$

Then ρ_ is minimax-optimal.*

Proof of Lemma 8.5. For any decision procedure ρ_o ,

$$\sup_{\theta \in \Theta} R(\rho_o, \theta) \geq \limsup_{k \rightarrow \infty} R(\rho_o, \Pi_k) \geq \lim_{k \rightarrow \infty} \inf_{\rho} R(\rho, \Pi_k) = \sup_{\theta \in \Theta} R(\rho_*, \theta).$$

\square

Lemma 8.6. *Let Θ be a topological space equipped with its Borel- σ -field. Suppose that any decision procedure ρ has continuous, real-valued risk function. If Π_o is a prior on Θ such that $\Pi_o(U) > 0$ for any non-void open set $U \subset \Theta$, and if ρ_o is a Bayes-optimal decision procedure for Π_o with $R(\rho_o, \Pi_o) < \infty$, then ρ_o is admissible.*

Exercise 8.7. Suppose that each distribution \mathbb{P}_{θ} has a density f_{θ} with respect to some σ -finite measure M on (Ω, \mathcal{A}) such that for every $\omega \in \Omega$, $\theta \rightarrow f_{\theta}(\omega)$ is continuous on Θ . Further, suppose that the loss function L is bounded, and for any $v \in \mathbf{V}$, let $\theta \mapsto L(v, \theta)$ be continuous on Θ . Show that each decision procedure has bounded and continuous risk function. (Hint: Use Scheffé's theorem and dominated convergence.)

Proof of Lemma 8.6. Suppose that ρ_o is not admissible. That means, there exists a decision procedure ρ such that $R(\rho, \cdot) \leq R(\rho_o, \cdot)$ and $R(\rho, \theta_o) < R(\rho_o, \theta_o)$ for some $\theta_o \in \Theta$. Since both risk functions $R(\rho, \cdot)$ and $R(\rho_o, \cdot)$ are continuous, there exist an open set $U \subset \Theta$ and an $\epsilon > 0$ such that $R(\rho, \cdot) \leq R(\rho_o, \cdot) - \epsilon$ on U . But then

$$\begin{aligned} R(\rho, \Pi_o) &= \int_{\Theta \setminus U} R(\rho, \theta) \Pi(d\theta) + \int_U R(\rho, \theta) \Pi(d\theta) \\ &\leq \int_{\Theta \setminus U} R(\rho_o, \theta) \Pi(d\theta) + \int_U (R(\rho_o, \theta) - \epsilon) \Pi(d\theta) \\ &= R(\rho_o, \Pi_o) - \epsilon \Pi_o(U) \\ &< R(\rho_o, \Pi_o), \end{aligned}$$

a contradiction to Bayes-optimality of ρ_o . \square

Example 8.8 (Gaussian location model). For a given sample size $n \in \mathbb{N}$ and a given standard deviation $\sigma > 0$, let $\mathcal{E} = (\mathcal{N}(\mu, \sigma^2)^{\otimes n})_{\mu \in \mathbb{R}}$. The sample mean $\hat{\mu}_*(\mathbf{x}) := \bar{x}$ is a minimax-optimal point estimator of $g(\mu) = \mu$, if

$$L(v, \mu) := (v - \mu)^2.$$

To see this, note first that $R(\hat{\mu}_*, \cdot) \equiv \sigma^2/n$. By Exercise 8.3, it suffices to consider non-randomized point estimators. Moreover, if $\Pi_k = \mathcal{N}(0, k)$, we know from Example 6.20 that the Bayes-optimal estimator is given by $\hat{\mu}_k(\mathbf{x}) = n\bar{x}/(n + \sigma^2/k)$ with Bayes-risk $R(\hat{\mu}_k, \Pi_k) = \sigma^2/(n + \sigma^2/k)$. Since this converges to σ^2/n as $k \rightarrow \infty$, Lemma 8.5 shows that $\hat{\mu}_*$ is minimax-optimal.

Exercise 8.9 (Point estimation of a binomial parameter). For $p \in [0, 1]$ let $\mathbb{P}_p := \text{Bin}(n, p)$. We consider point estimators of p with loss function $L : [0, 1] \times [0, 1] \rightarrow [0, \infty)$ given by

$$L(v, p) := (v - p)^2.$$

(a) Consider the Bayesian model with a random parameter $p \sim \text{Beta}(a, b)$ with given “hyper-parameters” $a, b > 0$, and a random observation X with $\mathcal{L}(X | p) = \mathbb{P}_p$. Here $\text{Beta}(a, b)$ is the distribution on $(0, 1)$ with Lebesgue density

$$\beta_{a,b}(p) = \frac{p^{a-1}(1-p)^{b-1}}{B(a, b)}, \quad B(a, b) := \int_0^1 u^{a-1}(1-u)^{b-1} du = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

Show that

$$\mathcal{L}(p | X) = \text{Beta}(a + X, b + n - X) \quad \text{and} \quad \mathbb{E}(p | X) = \frac{a + X}{a + b + n}.$$

(b) Adapt arguments from previous chapters to show that the estimator

$$\hat{p}_{a,b}(x) := \frac{a + x}{a + b + n}$$

minimizes the Bayes risk

$$R(\hat{p}, \text{Beta}(a, b)) := \int_0^1 R(\hat{p}, p) \text{Beta}(a, b)(dp).$$

(c) Determine the risk function $R(\hat{p}_{a,b}, \cdot)$ of the estimator $\hat{p}_{a,b}$ in part (b) explicitly.

(d) Now find parameters a, b such that the risk function in part (c) is constant. What are the consequences for the corresponding prior $\text{Beta}(a, b)$ and the corresponding estimator $\hat{p}_{a,b}$?

Unbiasedness. Sometimes it is difficult to find decision procedures satisfying some optimality criterion. But often the problem gets easier if we impose some additional constraints which are quite natural by themselves. Here is one such constraint: A decision procedure ρ is called *unbiased*, if for arbitrary $\theta, \eta \in \Theta$,

$$(8.2) \quad R(\rho, \theta) = \int_{\Omega} \int_{\mathbb{V}} L(v, \theta) \rho(\omega, dv) \mathbb{P}_{\theta}(d\omega) \leq \int_{\Omega} \int_{\mathbb{V}} L(v, \eta) \rho(\omega, dv) \mathbb{P}_{\theta}(d\omega).$$

Example 8.1 (Point estimation, continued). Let $L(v, \theta) := \|v - g(\theta)\|^2$ with some Euclidean norm $\|\cdot\|$ on \mathbb{R}^q , that is, $\|x\| = \sqrt{x^{\top}Ax}$ for some symmetric, positive definite matrix $A \in \mathbb{R}^{q \times q}$.

For any estimator $\hat{g} : \Omega \rightarrow \mathbb{R}^q$, the risk $R(\hat{g}, \theta) = \mathbb{E}_\theta(\|\hat{g} - g(\theta)\|^2)$ is finite if, and only if, $\mathbb{E}_\theta(\|\hat{g}\|^2) < \infty$. In the latter case,

$$\begin{aligned}\mathbb{E}_\theta(\|\hat{g} - g(\eta)\|^2) &= \mathbb{E}_\theta(\|\hat{g} - \mathbb{E}_\theta(\hat{g})\|^2) + \|\mathbb{E}_\theta(\hat{g}) - g(\eta)\|^2, \\ \mathbb{E}_\theta(\|\hat{g} - g(\theta)\|^2) &= \mathbb{E}_\theta(\|\hat{g} - \mathbb{E}_\theta(\hat{g})\|^2) + \|\mathbb{E}_\theta(\hat{g}) - g(\theta)\|^2,\end{aligned}$$

so \hat{g} is unbiased if and only if for all $\theta \in \Theta$,

$$\|\mathbb{E}_\theta(\hat{g}) - g(\theta)\| = \min_{\eta \in \Theta} \|\mathbb{E}_\theta(\hat{g}) - g(\eta)\|.$$

Consequently, if \hat{g} is an estimator such that

$$\mathbb{E}_\theta(\|\hat{g}\|^2) < \infty \quad \text{and} \quad \mathbb{E}_\theta(\hat{g}) \in \text{closure}\{g(\eta) : \eta \in \Theta\}$$

for all $\theta \in \Theta$, then \hat{g} is unbiased if, and only if,

$$\mathbb{E}_\theta(\hat{g}) = g(\theta) \quad \text{for all } \theta \in \Theta.$$

Example 8.2 (Statistical tests, continued). We consider the general loss function $L(v, \theta) = 1_{[v=1, \theta \in \Theta_0]} \lambda_1 + 1_{[v=0, \theta \in \Theta_1]} \lambda_2$ with $\lambda_1, \lambda_2 > 0$. Then

$$\int_{\Omega} \int_{\mathbb{V}} L(v, \eta) \rho(\omega, dv) \mathbb{P}_\theta(d\omega) \geq R(\varphi, \theta) \quad \text{for arbitrary } \theta, \eta \in \Theta$$

is easily shown to be equivalent to

$$\begin{cases} \lambda_2(1 - \mathbb{E}_{\theta_0}(\varphi)) \geq \lambda_1 \mathbb{E}_{\theta_0}(\varphi) & \text{for any } \theta_0 \in \Theta_0, \\ \lambda_1 \mathbb{E}_{\theta_1}(\varphi) \geq \lambda_2(1 - \mathbb{E}_{\theta_1}(\varphi)) & \text{for any } \theta_1 \in \Theta_1. \end{cases}$$

In other words,

$$\mathbb{E}_\theta(\varphi) \begin{cases} \leq \alpha & \text{for } \theta \in \Theta_0 \\ \geq \alpha & \text{for } \theta \in \Theta_1 \end{cases}$$

with

$$\alpha := \frac{\lambda_2}{\lambda_1 + \lambda_2} \in (0, 1).$$

Coming from the other end, for a given test level $\alpha \in (0, 1)$, a level- α test of Θ_0 with power at least α for each $\theta \in \Theta_A$ is unbiased in the sense of (8.2) if, say, $\lambda_1 = 1 - \alpha$ and $\lambda_2 = \alpha$.

Exercise 8.10 (Estimating functions of a binomial parameter). Unbiasedness of point estimators seems often a natural constraint. But it is potentially too restrictive. Consider the statistical experiment $(\text{Bin}(n, p))_{p \in [0, 1]}$ and an arbitrary function $g : [0, 1] \rightarrow \mathbb{R}$.

(a) Suppose that $\hat{g} : \{0, 1, \dots, n\} \rightarrow \mathbb{R}$ is an unbiased estimator of $g(p)$, i.e. $\mathbb{E}_p(\hat{g}) = g(p)$ for all $p \in [0, 1]$. Show that $g(p)$ is a polynomial in p of order at most n .

(b) Show that an unbiased estimator \hat{g} as in part (a) is unique.

(c) Determine the estimator \hat{g} explicitly in case of $g(p) = p^k$ for some $k \in \{1, \dots, n\}$. Hint: Consider factorials $[x]_k := \prod_{0 \leq i < k} (x - i)$.

8.3 Informativity and Sufficiency

This section is about the comparison of two statistical experiments with one and the same parameter space Θ .

8.3.1 Informativity

Let $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ and $\tilde{\mathcal{E}} = (\tilde{\Omega}, \tilde{\mathcal{A}}, (\tilde{\mathbb{P}}_\theta)_{\theta \in \Theta})$ be statistical experiments. Experiment $\tilde{\mathcal{E}}$ is called *at least as informative as* experiment \mathcal{E} , if the following condition is satisfied: Let $(\mathbb{V}, \mathcal{B})$ be an arbitrary decision space and $L : \mathbb{V} \times \Theta \rightarrow (-\infty, \infty]$ be any loss function. Then for any decision procedure $\rho : \Omega \times \mathcal{B} \rightarrow [0, 1]$, there exists a decision procedure $\tilde{\rho} : \tilde{\Omega} \times \mathcal{B} \rightarrow [0, 1]$ such that

$$R(\tilde{\rho}, \cdot) \leq R(\rho, \cdot) \quad \text{on } \Theta.$$

This looks like a very strong condition, because it involves arbitrary decision spaces and loss functions. Nevertheless there is an elegant criterion due to David Blackwell:

Lemma 8.11 (Blackwell's criterion). *Let $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ and $\tilde{\mathcal{E}} = (\tilde{\Omega}, \tilde{\mathcal{A}}, (\tilde{\mathbb{P}}_\theta)_{\theta \in \Theta})$ be statistical experiments. Suppose there exists a stochastic kernel K from $(\tilde{\Omega}, \tilde{\mathcal{A}})$ to (Ω, \mathcal{A}) such that for any $\theta \in \Theta$ and $A \in \mathcal{A}$,*

$$\mathbb{P}_\theta(A) = \tilde{\mathbb{P}}_\theta \otimes K(\tilde{\Omega} \times A) = \int_{\tilde{\Omega}} K(\tilde{\omega}, A) \tilde{\mathbb{P}}_\theta(d\tilde{\omega}).$$

Then $\tilde{\mathcal{E}}$ is at least as informative as \mathcal{E} .

The intuition behind Blackwell's criterion is as follows: Suppose that Bob and Alice are planning statistical experiments \mathcal{E} and $\tilde{\mathcal{E}}$, respectively, to do inference about an unknown parameter in Θ .

Alice: "My experiment is at least as good as yours!" Bob: "How do you know?"

Alice: "Well, tell me any decision space and loss function." Bob: " $(\mathbb{V}, \mathcal{B})$ and L ."

Alice: "Okay, and what is your favourite decision procedure for that?" Bob: " ρ "

Alice: "All right, this is what I will do: If my experiment yields data $\tilde{X} \sim \tilde{\mathbb{P}}_\theta$, where $\theta \in \Theta$ is unknown, I will use Mr. Blackwell's kernel K to generate new data $X \sim K(\tilde{X}, \cdot)$. And then, I'll use your ρ to draw a decision $V \sim \rho(X, \cdot)$!" Bob: "That's cheap!"

Proof of Lemma 8.11. The assumption implies that for any measurable function $h : (\Omega, \mathcal{A}) \rightarrow [0, \infty]$,

$$\int_{\Omega} h d\mathbb{P}_\theta = \int_{\tilde{\Omega}} \int_{\Omega} h(\omega) K(\tilde{\omega}, d\omega) \tilde{\mathbb{P}}_\theta(d\tilde{\omega}).$$

By assumption, this is true for indicator functions $h = 1_A$ with $A \in \mathcal{A}$. By linearity of integration, this is true for measurable functions taking only finitely many different values in $[0, \infty)$. By monotone convergence, the asserted equation is true for arbitrary measurable functions $h : (\Omega, \mathcal{A}) \rightarrow [0, \infty]$.

For a given decision space $(\mathbb{V}, \mathcal{B})$ and decision procedure $\rho : \Omega \times \mathcal{B} \rightarrow [0, 1]$ we define a decision procedure $\tilde{\rho} : \tilde{\Omega} \times \mathcal{B} \rightarrow [0, 1]$ as follows: For $B \in \mathcal{B}$ we set

$$\tilde{\rho}(\tilde{\omega}, B) := \int_{\Omega} \rho(\omega, B) K(\tilde{\omega}, d\omega).$$

This construction of $\tilde{\rho}$ implies that for any measurable function $h : (\mathbb{V}, \mathcal{B}) \rightarrow [0, \infty]$ and $\tilde{\omega} \in \tilde{\Omega}$,

$$\int_{\mathbb{V}} h(v) \tilde{\rho}(\tilde{\omega}, dv) = \int_{\Omega} \int_{\mathbb{V}} h(v) \rho(\omega, dv) K(\tilde{\omega}, d\omega).$$

Hence for loss functions $L : \mathbb{V} \times \Theta \rightarrow [0, \infty]$ and any $\theta \in \Theta$,

$$\begin{aligned} R(\tilde{\rho}, \theta) &= \int_{\tilde{\Omega}} \int_{\mathbb{V}} L(v, \theta) \tilde{\rho}(\tilde{\omega}, dv) \tilde{\mathbb{P}}_{\theta}(d\tilde{\omega}) \\ &= \int_{\tilde{\Omega}} \int_{\Omega} \int_{\mathbb{V}} L(v, \theta) \rho(\omega, dv) K(\tilde{\omega}, d\omega) \tilde{\mathbb{P}}_{\theta}(d\tilde{\omega}) \\ &= \int_{\Omega} \int_{\mathbb{V}} L(v, \theta) \rho(\omega, dv) \mathbb{P}_{\theta}(d\omega) \\ &= R(\rho, \theta). \end{aligned}$$

The same equation is true for loss functions L with values in $(-\infty, \infty]$, provided that $R(\rho', \theta)$ or $R(\rho, \theta)$ is well-defined in $(-\infty, \infty]$. Just write $L = L^+ - L^-$ with $L^{\pm} := \max(\pm L, 0)$. \square

Example 8.12 (Sampling with and without replacement). Let \mathcal{M} be a population of *known* size $N = \#\mathcal{M}$, but with unknown characteristics $\theta \in \Theta$, the latter describing certain properties of the individuals in \mathcal{M} . Suppose we draw a random sample of size $n \geq 2$ from \mathcal{M} *with* replacement. That means, we obtain a sample $\omega = (\omega_1, \dots, \omega_n)$ in the set $\Omega := \mathcal{M}^n$ with N^n different elements. Here the corresponding distributions \mathbb{P}_{θ} , $\theta \in \Theta$, are all identical: the uniform distribution on Ω . The subscript θ just indicates that the potential samples may differ in terms of certain characteristics of the individuals.

In case of $N \geq n$, sampling *without* replacement would be an alternative strategy, leading to the experiment $\tilde{\mathcal{E}} = (\tilde{\Omega}, \mathcal{P}(\tilde{\Omega}), (\tilde{\mathbb{P}}_{\theta})_{\theta \in \Theta})$ with $\tilde{\mathbb{P}}_{\theta}$ denoting the uniform distribution on

$$\tilde{\Omega} := \{\tilde{\omega} \in \Omega : \tilde{\omega}_i \neq \tilde{\omega}_j \text{ for } 1 \leq i < j \leq n\},$$

a set of $[N]_n = N(N-1) \cdots (N-n+1)$ different samples. Indeed the experiment $\tilde{\mathcal{E}}$ is at least as informative as $\mathcal{E} = (\Omega, \mathcal{P}(\Omega), (\mathbb{P}_{\theta})_{\theta \in \Theta})$. For if we draw a sample $\tilde{\omega} \in \tilde{\Omega}$ from the uniform distribution on $\tilde{\Omega}$, we may generate a new sample $\omega \in \Omega$ as follows: We choose $\omega_1 := \tilde{\omega}_1$ and independent random points $\omega_2, \dots, \omega_n$ with distribution

$$\omega_j \sim \sum_{i=1}^{j-1} \frac{1}{N} \delta_{\tilde{\omega}_i} + \frac{N-j+1}{N} \delta_{\tilde{\omega}_j}.$$

In other words, we go through the “list” $\tilde{\omega}$, and each component $\tilde{\omega}_j$ is kept with probability $1 - (j-1)/N$ or replaced with a randomly chosen predecessor $\tilde{\omega}_i$, $i < j$. This defines a probability distribution $K(\tilde{\omega}, \cdot)$ on Ω such that the resulting ω is uniformly distributed on Ω .

Exercise 8.13 (From sampling without to sampling with replacement). Show that the construction of ω from $\tilde{\omega}$ in Example 8.12 corresponds to a stochastic kernel K from $\tilde{\Omega}$ to Ω such that

$$K(\tilde{\omega}, \{\tilde{\omega}_1, \dots, \tilde{\omega}_n\}^n) = 1$$

and for any $y \in \mathcal{M}^n$,

$$\frac{1}{[N]^n} \sum_{\tilde{\omega} \in \tilde{\Omega}} K(\tilde{\omega}, \{y\}) = \frac{1}{N^n}.$$

Exercise 8.14 (Estimating the reciprocal of a population size). At first glance one could think that sampling from a population \mathcal{M} without replacement is always more informative than sampling with replacement. But this is not true in general. For instance, suppose that the size N of the population \mathcal{M} is unknown. Then drawing a random sample of size $n \leq N$ without replacement from that population reveals nothing about the population size N , but sampling with replacement yields some information:

For a given integer $n \geq 2$, let $\mathbb{P}_{\mathcal{M}}$ be the uniform distribution on \mathcal{M}^n . We are interested in constructing an unbiased estimator of $g(\mathcal{M}) := 1/N$.

(a) Determine the expectation of X_i where

$$X_i(\omega) := \frac{1_{[\omega_i \in \{\omega_j : j \neq i\}]}}{\#\{\omega_1, \omega_2, \dots, \omega_n\}}, \quad \text{for } \omega \in \mathcal{M}^n.$$

for all $\omega \in \mathcal{M}^n$ and $i \in \{1, 2, \dots, n\}$.

(b) Propose an unbiased estimator of $g(\mathcal{M})$.

(c) For $1 \leq i < j \leq n$, let $X_{ij}(\omega) := 1_{[\omega_i = \omega_j]}$. Determine the expectation of X_{ij} and propose an unbiased estimator of $g(\mathcal{M})$.

(d) Determine the standard deviation of your estimator in part (c).

8.3.2 Sufficiency

The concept of sufficiency is a special instance of Blackwell's criterion. We consider a statistical experiment $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_{\theta})_{\theta \in \Theta})$. Now we want to know whether it is sufficient to restrict one's attention to partial information about the experiment's outcome. Partial information could mean that we want to replace an observation $X \sim \mathbb{P}_{\theta}$ with a given function $T(X)$, or we want to restrict our attention to events A_o in a sub- σ -field \mathcal{A}_o of \mathcal{A} .

Definition 8.15 (Sufficient statistic). A measurable mapping $T : (\Omega, \mathcal{A}) \rightarrow (\tilde{\Omega}, \tilde{\mathcal{A}})$ is called a *sufficient statistic for \mathcal{E}* , if there exists a stochastic kernel K from $(\tilde{\Omega}, \tilde{\mathcal{A}})$ to (Ω, \mathcal{A}) describing the conditional distribution of $X \sim \mathbb{P}_{\theta}$, given $T(X)$, for any $\theta \in \Theta$. In other words, for arbitrary $\theta \in \Theta$, $A \in \mathcal{A}$ and $\tilde{A} \in \tilde{\mathcal{A}}$,

$$\mathbb{P}_{\theta}(\{T \in \tilde{A}\} \cap A) = \int_{\tilde{A}} K(t, A) \mathbb{P}_{\theta}^T(dt).$$

Sufficiency of T implies that the experiment $\mathcal{E}^T := (\tilde{\Omega}, \tilde{\mathcal{A}}, (\mathbb{P}_{\theta}^T)_{\theta \in \Theta})$ is at least as informative as \mathcal{E} . In other words, one may reduce raw data $X \sim \mathbb{P}_{\theta}$ to $T(X)$ without any loss of information.

Indeed, having reduced X to $T(\mathbf{X})$, one could generate an artificial observation $\tilde{X} \sim K(T(X), \cdot)$, and $\tilde{X} \sim \mathbb{P}_\theta$, too, no matter what value the unknown parameter $\theta \in \Theta$ has. Sufficiency can often be verified with the following criterion:

Theorem 8.16 (Neyman's factorization criterion). *Suppose that (Ω, d) is a separable and complete metric space, and let $\mathcal{A} = \text{Borel}(\Omega, d)$. Further let M be a σ -finite measure on (Ω, \mathcal{A}) such that each distribution \mathbb{P}_θ has a density f_θ with respect to M . Suppose that there exists a measurable function $h : (\Omega, \mathcal{A}) \rightarrow [0, \infty)$ such that for any $\theta \in \Theta$ and $\omega \in \Omega$,*

$$f_\theta(\omega) = g_\theta(T(\omega))h(\omega)$$

with $g_\theta : (\tilde{\Omega}, \tilde{\mathcal{A}}) \rightarrow [0, \infty)$ measurable. Then T is a sufficient statistic for \mathcal{E} .

Corollary 8.17. *Let Ω and $\tilde{\Omega}$ be countable sets equipped with $\mathcal{A} = \mathcal{P}(\Omega)$ and $\tilde{\mathcal{A}} = \mathcal{P}(\tilde{\Omega})$. Suppose that there exists a function $h : \Omega \rightarrow [0, \infty)$ such that for any $\theta \in \Theta$ and $\omega \in \Omega$,*

$$\mathbb{P}_\theta(\{\omega\}) = g_\theta(T(\omega))h(\omega)$$

with $g_\theta : \tilde{\Omega} \rightarrow [0, \infty)$. Then T is a sufficient statistic for \mathcal{E} .

For readers feeling uneasy about measure theory it may be instructive to prove the latter corollary directly. It is a consequence of Theorem 8.16 if we use the metric $d(\tilde{\omega}, \tilde{\omega}') := 1_{[\tilde{\omega} \neq \tilde{\omega}]}$ on $\tilde{\Omega}$ and the counting measure M on Ω , that is, $M(\{\omega\}) = 1$ for all $\omega \in \Omega$.

Proof of Theorem 8.16. A measure M on (Ω, \mathcal{A}) is σ -finite if, and only if, there exists a measurable function $J : (\Omega, \mathcal{A}) \rightarrow (0, \infty)$ such that $\int J dM = 1$. But then we could replace M with the measure $J \cdot M$, i.e. $A \mapsto \int_A J dM$, and h with h/J . Hence we may and do assume without loss of generality that M is a probability measure. By our assumption on (Ω, \mathcal{A}) there exists a stochastic kernel K_o from $(\tilde{\Omega}, \tilde{\mathcal{A}})$ to (Ω, \mathcal{A}) such that for arbitrary $\tilde{A} \in \tilde{\mathcal{A}}$ and $A \in \mathcal{A}$,

$$M(\{T \in \tilde{A}\} \cap A) = \int_{\tilde{A}} K_o(t, A) M^T(dt),$$

see Chapter 4. More generally, for arbitrary measurable and non-negative functions f on $(\tilde{\Omega} \times \Omega, \tilde{\mathcal{A}} \otimes \mathcal{A})$,

$$\int_{\Omega} f(T(\omega), \omega) M(d\omega) = \int_{\tilde{\Omega}} \int_{\Omega} f(t, \omega) K_o(t, d\omega) M^T(dt).$$

This implies that for arbitrary $\theta \in \Theta$, $\tilde{A} \in \tilde{\mathcal{A}}$ and $A \in \mathcal{A}$,

$$\begin{aligned} \mathbb{P}_\theta(\{T \in \tilde{A}\} \cap A) &= \int_{\Omega} 1_{\tilde{A}}(T(\omega)) 1_A(\omega) g_\theta(T(\omega)) h(\omega) M(d\omega) \\ &= \int_{\tilde{\Omega}} \int_{\Omega} 1_{\tilde{A}}(t) g_\theta(t) 1_A(\omega) h(\omega) K_o(t, d\omega) M^T(dt) \\ &= \int_{\tilde{A}} g_\theta(t) \int_A h(\omega) K_o(t, d\omega) M^T(dt). \end{aligned}$$

Taking $A = \Omega$ shows that the density of \mathbb{P}_θ^T with respect to M^T is given by

$$f_\theta^T(t) := g_\theta(t)H(t) \quad \text{with} \quad H(t) := \int_{\Omega} h(\omega) K_o(t, d\omega).$$

In particular, the set $N := \{t \in \tilde{\Omega} : H(t) = 0 \text{ or } H(t) = \infty\}$ satisfies $\mathbb{P}_\theta^T(N) = 0$ for any $\theta \in \Theta$, because $\{H = 0\} \subset \{f_\theta^T = 0\}$ and $\{H = \infty\} \subset \{f_\theta^T = 0\} \cup \{f_\theta^T = \infty\}$, and both $\mathbb{P}_\theta^T(f_\theta^T = 0)$ and $\mathbb{P}_\theta^T(f_\theta^T = \infty)$ are equal to zero. Hence

$$K(t, A) := \begin{cases} H(t)^{-1} \int_A h(\omega) K_o(t, d\omega) & \text{if } 0 < H(t) < \infty \\ M(A) & \text{else} \end{cases}$$

defines a stochastic kernel K from $(\tilde{\Omega}, \tilde{\mathcal{A}})$ to (Ω, \mathcal{A}) such that

$$\mathbb{P}_\theta(\{T \in \tilde{A}\} \cap A) = \int_{\tilde{A}} f_\theta^T(t) K(t, A) M^T(dt) = \int_{\tilde{A}} K(t, A) \mathbb{P}_\theta^T(dt)$$

for arbitrary $\theta \in \Theta$, $\tilde{A} \in \tilde{\mathcal{A}}$ and $A \in \mathcal{A}$. □

Example 8.18 (Bernoulli sequences). Let X_1, X_2, \dots, X_n be independent, identically distributed random variables with values in $\{0, 1\}$ and unknown parameter $p = \mathbb{P}(X_i = 1) = \mathbb{E}(X_i) \in [0, 1]$. This leads to the statistical experiment $\mathcal{E} = (\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\mathbb{P}_p)_{p \in [0, 1]})$ with $\mathbb{P}_p := ((1-p)\delta_0 + p\delta_1)^{\otimes n}$. The density f_p of \mathbb{P}_p with respect to counting measure on $\{0, 1\}^n$ is given by

$$f_p(\omega) = \prod_{i=1}^n p^{\omega_i} (1-p)^{1-\omega_i} = p^{T(\omega)} (1-p)^{n-T(\omega)},$$

i.e. a function of $T(\omega) := \sum_{i=1}^n \omega_i$ only. Hence T is a sufficient statistic for \mathcal{E} . Indeed, T has distribution $\text{Bin}(n, p)$, and for any $t \in \{0, 1, \dots, n\}$, the conditional distribution $\mathbb{P}_p(\cdot | T = t)$ is the uniform distribution on the set of all $\omega \in \{0, 1\}^n$ with $T(\omega) = t$.

Example 8.19 (Gaussian samples). Let X_1, X_2, \dots, X_n be independent, identically distributed random variables with distribution $\mathcal{N}(\mu, \sigma^2)$, the mean $\mu \in \mathbb{R}$ and the standard deviation $\sigma > 0$ being unknown. This leads to the statistical experiment $\mathcal{E} = (\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (\mathbb{P}_\theta)_{\theta \in \Theta})$ with $\Theta = \mathbb{R} \times (0, \infty)$ and $\mathbb{P}_{\mu, \sigma} := \mathcal{N}(\mu, \sigma^2)^{\otimes n}$. Setting $X_i(\omega) := \omega_i$, the density $f_{\mu, \sigma}$ of $\mathbb{P}_{\mu, \sigma}$ with respect to Lebesgue measure on \mathbb{R}^n times $(2\pi)^{-n/2}$ is given by

$$\begin{aligned} f_{\mu, \sigma}(\mathbf{x}) &= \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} - n \log \sigma\right) \\ &= \exp\left(-\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{2\sigma^2} - \frac{n(\bar{x} - \mu)^2}{2\sigma^2} - n \log \sigma\right) \\ &= g_{\mu, \sigma}(T(\mathbf{x})), \end{aligned}$$

where $T(\mathbf{x}) = (T_1(\mathbf{x}), T_2(\mathbf{x}))$ and

$$\begin{aligned} T_1(\mathbf{x}) &:= \bar{x} = n^{-1} \sum_{i=1}^n x_i, \\ T_2(\mathbf{x}) &:= \sum_{i=1}^n (x_i - \bar{x})^2, \\ g_{\mu, \sigma}(t_1, t_2) &:= \exp\left(-\frac{n(t_1 - \mu)^2}{2\sigma^2} - \frac{t_2}{2\sigma^2} - n \log \sigma\right). \end{aligned}$$

Hence the statistic $T = (T_1, T_2) : \mathbb{R}^n \rightarrow \mathbb{R} \times [0, \infty)$ is sufficient for the experiment \mathcal{E} .

It is worthwhile here to verify sufficiency directly, based on standard arguments in connection with student's t distribution: Let $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n$ be an orthonormal basis of \mathbb{R}^n , where $\mathbf{b}_1 = n^{-1/2}\mathbf{1}_n$. Then the distribution $\mathbb{P}_{\mu, \sigma}$ coincides with the distribution of

$$\tilde{\mathbf{X}} := \mu\mathbf{1}_n + \sigma \sum_{i=1}^n Z_i \mathbf{b}_i$$

with stochastically independent random variables $Z_1, \dots, Z_n \sim \mathcal{N}(0, 1)$. The statistic $T(\tilde{\mathbf{X}})$ is equal to

$$\left(\mu + n^{-1/2}\sigma Z_1, \sigma^2 \sum_{i=2}^n Z_i^2 \right),$$

and we may write

$$\tilde{\mathbf{X}} = T_1(\tilde{\mathbf{X}})\mathbf{1}_n + \sqrt{T_2(\tilde{\mathbf{X}})} \sum_{i=2}^n W_i \mathbf{b}_i$$

with

$$W_i := \left(\sum_{j=2}^n Z_j^2 \right)^{-1/2} Z_i.$$

The random vector $\mathbf{W} := (W_i)_{i=2}^n$ is uniformly distributed on the unit sphere in \mathbb{R}^{n-1} and stochastically independent from $\sum_{i=2}^n Z_i^2$. Hence the conditional distribution of $\tilde{\mathbf{X}}$, given \tilde{T} , does not depend on the parameter (μ, σ) .

Exercise 8.20 (Gamma distributions). Let $\mathbf{X} = (X_i)_{i=1}^n$ have $n \geq 2$ independent components with

$$X_i \sim \text{Gamma}(a, b)$$

and unknown parameters $a, b > 0$.

(a) Determine an \mathbb{R}^2 -valued sufficient statistic $T(\mathbf{X})$ for the corresponding statistical experiment $(\text{Gamma}(a, b))_{a, b > 0}$.

(b) Determine the conditional distribution of \mathbf{X} given $T(\mathbf{X})$ in case of $n = 2$.

Exercise 8.21 (Markov chains with finite state space). Let $\mathbf{X} = (X_t)_{t=0}^n$ be a Markov chain with values in a finite set \mathcal{X} and fixed starting point $X_0 = x_0 \in \mathcal{X}$. That means, for $1 \leq k < n$ and $y_0, \dots, y_n, z \in \mathcal{X}$

$$\mathbb{P}_\theta(X_{k+1} = z \mid (X_t)_{t=0}^k = (y_t)_{t=0}^k) = \theta_{y_k, z}$$

with an unknown “matrix” $\theta \in [0, 1]^{\mathcal{X} \times \mathcal{X}}$ such that

$$\sum_{z \in \mathcal{X}} \theta_{y, z} = 1 \quad \text{for all } y \in \mathcal{X}.$$

Let Θ be the set of all such “matrices” θ . Determine a sufficient statistic for the statistical experiment $(\Omega, \mathcal{P}(\Omega), (\mathbb{P}_\theta)_{\theta \in \Theta})$, where $\Omega = \{\mathbf{y} \in \mathcal{X}^{\{0, 1, \dots, n\}} : y_0 = x_0\}$.

Exercise 8.22. Let \mathcal{M} be a population of individuals with identification numbers in \mathbb{Z} . We assume that the set of all identification numbers equals $\{a, \dots, b\}$ with unknown integers $a \leq b$. We only know that $b - a > n$ for some given integer $n \geq 2$.

Now we draw a random sample of size n without replacement from \mathcal{M} and note the tuple $\omega = (\omega_1, \omega_2, \dots, \omega_n)$ of identification numbers.

Show that $T(\omega) = (\min(\omega), \max(\omega))$ is a sufficient statistic for this experiment. Describe the conditional distribution of ω given $T(\omega) = (s_1, s_2)$.

Hint: The formal definition of the experiment is

$$\begin{aligned}\Omega &= \{\omega \in \mathbb{Z}^n : \omega_i \neq \omega_j \text{ whenever } i \neq j\}, \\ \mathbb{P}_{a,b} &= \text{Unif}(\Omega \cap \{a, \dots, b\}^n), \\ \Theta &= \{(a, b) \in \mathbb{Z}^n : b - a > n\}.\end{aligned}$$

Definition 8.23 (Sufficient sub- σ -fields). Let \mathcal{A}_o be a σ -field over Ω such that $\mathcal{A}_o \subset \mathcal{A}$. It is called a *sufficient sub- σ -field* for \mathcal{E} , if there exists a stochastic kernel K from (Ω, \mathcal{A}_o) to (Ω, \mathcal{A}) describing the conditional distribution of $X \sim \mathbb{P}_\theta$, given \mathcal{A}_o , for any $\theta \in \Theta$. In other words, for arbitrary $\theta \in \Theta$, $A \in \mathcal{A}$ and $A_o \in \mathcal{A}_o$,

$$\mathbb{P}_\theta(A_o \cap A) = \int_{A_o} K(\omega, A) \mathbb{P}_\theta(d\omega).$$

Note that sufficiency of \mathcal{A}_o is equivalent to sufficiency of the statistic

$$T : (\Omega, \mathcal{A}) \rightarrow (\Omega, \mathcal{A}_o), \quad T(\omega) := \omega.$$

Sufficiency of \mathcal{A}_o implies that the experiment $\mathcal{E}_o := (\Omega, \mathcal{A}_o, (\mathbb{P}_\theta)_{\theta \in \Theta})$ is at least as informative as \mathcal{E} . In other words, when analyzing raw data $X \sim \mathbb{P}_\theta$, we may restrict our attention to decision procedures which are \mathcal{A}_o -measurable rather than \mathcal{A} -measurable.

Example 8.24 (Invariant distributions). Let \mathcal{G} be a finite group of measurable bijective mappings $g : (\Omega, \mathcal{A}) \rightarrow (\Omega, \mathcal{A})$. That means, for arbitrary $g, h \in \mathcal{G}$, both $h \circ g$ and g^{-1} belong to \mathcal{G} , too. Now let $\mathcal{A}_\mathcal{G}$ be the set of \mathcal{G} -invariant sets $A \in \mathcal{A}$, i.e.

$$g(A) = \{g(\omega) : \omega \in A\} = A \quad \text{for all } g \in \mathcal{G}.$$

This is obviously a sub- σ -field of \mathcal{A} .

Now suppose that all distributions \mathbb{P}_θ , $\theta \in \Theta$, are \mathcal{G} -invariant in the sense that

$$\mathbb{P}_\theta^g = \mathbb{P}_\theta \quad \text{for all } g \in \mathcal{G},$$

where \mathbb{P}_θ^g is the image measure $\mathbb{P}_\theta \circ g^{-1}$. Then $\mathcal{A}_\mathcal{G}$ is sufficient for \mathcal{E} , and the conditional distribution $\mathbb{P}_\theta(\cdot | \mathcal{A}_o)$ is given by the stochastic kernel

$$K(\omega, A) := \frac{1}{\#\mathcal{G}} \sum_{g \in \mathcal{G}} 1_A(g(\omega))$$

i.e.

$$K(\omega, \cdot) = \frac{1}{\#\mathcal{G}} \sum_{g \in \mathcal{G}} \delta_{g(\omega)}.$$

Hence the experiment $\mathcal{E}_{\mathcal{G}} := (\Omega, \mathcal{A}_{\mathcal{G}}, (\mathbb{P}_{\theta})_{\theta \in \Theta})$ is at least as informative as \mathcal{E} .

Proof: Obviously, $K(\omega, \cdot)$ is a probability measure on (Ω, \mathcal{A}) for any $\omega \in \Omega$. For fixed $A \in \mathcal{A}$, the function $\omega \rightarrow K(\omega, A)$ is certainly \mathcal{A} -measurable. To verify $\mathcal{A}_{\mathcal{G}}$ -measurability it suffices to show that $K(h(\omega), A) = K(\omega, A)$ for arbitrary $\omega \in \Omega$ and $h \in \mathcal{G}$, see Exercise 8.25 below. But

$$K(h(\omega), A) = \frac{1}{\#\mathcal{G}} \sum_{g \in \mathcal{G}} 1_A(g \circ h(\omega)) = \frac{1}{\#\mathcal{G}} \sum_{\tilde{g} \in \mathcal{G}} 1_A(\tilde{g}(\omega)),$$

because for any fixed $h \in \mathcal{G}$, the mapping $\mathcal{G} \ni g \mapsto g \circ h \in \mathcal{G}$ is bijective, see Exercise 5.1. It remains to be shown that for arbitrary $\theta \in \Theta$, $A_o \in \mathcal{A}_{\mathcal{G}}$ and $A \in \mathcal{A}$,

$$\mathbb{P}_{\theta}(A_o \cap A) = \int_{A_o} K(\omega, A) \mathbb{P}_{\theta}(d\omega).$$

But the right hand side equals

$$\frac{1}{\#\mathcal{G}} \sum_{g \in \mathcal{G}} \int_{A_o} 1_A(g(\omega)) \mathbb{P}_{\theta}(d\omega),$$

and each summand equals

$$\begin{aligned} \int_{A_o} 1_A(g(\omega)) \mathbb{P}_{\theta}(d\omega) &= \mathbb{P}_{\theta}(A_o \cap g^{-1}(A)) \\ &= \mathbb{P}_{\theta}(g^{-1}(g(A_o) \cap A)) \\ &= \mathbb{P}_{\theta}(g^{-1}(A_o \cap A)) && \text{(by } \mathcal{G}\text{-invariance of } A_o\text{)} \\ &= \mathbb{P}_{\theta}^g(A_o \cap A) \\ &= \mathbb{P}_{\theta}(A_o \cap A) && \text{(by } \mathcal{G}\text{-invariance of } \mathbb{P}_{\theta}\text{)}. \end{aligned}$$

Exercise 8.25. In the setting of Example 8.24, let $(\mathbb{V}, \mathcal{B})$ be another measurable space, and let $\rho : (\Omega, \mathcal{A}) \rightarrow (\mathbb{V}, \mathcal{B})$ be a measurable mapping.

(a) Suppose that ρ is \mathcal{G} -invariant in the sense that $\rho \circ g \equiv \rho$ for arbitrary $g \in \mathcal{G}$. Show that ρ is $\mathcal{A}_{\mathcal{G}}$ - \mathcal{B} -measurable.

(a) Suppose that ρ is $\mathcal{A}_{\mathcal{G}}$ measurable, and suppose that \mathcal{B} “separates points in \mathbb{V} ”. That means, for arbitrary different points $v_1, v_2 \in \mathbb{V}$ there exists a set $B \in \mathcal{B}$ such that $v_1 \in B$ but $v_2 \notin B$. Show that ρ is \mathcal{G} -invariant.

Example 8.26 (Permutation-invariance). For some integer $n \geq 2$ let $(\Omega, \mathcal{A}) = (\mathcal{X}^n, \mathcal{B}^{\otimes n})$. Further let \mathcal{S}_n be the group of all permutations of $\{1, \dots, n\}$, that means, bijective mappings $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$. Any $\pi \in \mathcal{S}_n$ induces a measurable bijection $g_{\pi} : \Omega \rightarrow \Omega$,

$$\omega = (\omega_i)_{i=1}^n \mapsto g_{\pi}(\omega) := (\omega_{\pi(i)})_{i=1}^n.$$

Indeed one can easily show that for permutations $\pi, \sigma \in \mathcal{S}_n$,

$$g_{\pi} \circ g_{\sigma} = g_{\sigma \circ \pi}.$$

(Note also that $\pi \in \mathcal{S}_n$ is uniquely determined by the mapping g_π , unless $\#\mathcal{X} = 1$.) Thus $\mathcal{G} := \{g_\pi : \pi \in \mathcal{S}_n\}$ is a group of measurable bijections of (Ω, \mathcal{A}) .

A distribution \mathbb{P} on (Ω, \mathcal{A}) is called *permutation-invariant* or *exchangeable* if it is \mathcal{G} -invariant. In other words, if $X = (X_i)_{i=1}^n$ has distribution \mathbb{P} , then for any $\pi \in \mathcal{S}_n$, the random tuple $(X_{\pi(i)})_{i=1}^n$ has distribution \mathbb{P} , too.

If $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ consists of permutation-invariant distributions, then the sub- σ -field $\mathcal{A}_{\mathcal{G}}$ of all permutation-invariant sets is sufficient for \mathcal{E} . The corresponding stochastic kernel K from $(\Omega, \mathcal{A}_{\mathcal{G}})$ to (Ω, \mathcal{A}) is given by

$$K(\omega, \cdot) = \frac{1}{n!} \sum_{\pi \in \mathcal{S}_n} \delta_{g_\pi(\omega)}.$$

Example 8.27 (Sign invariance). For some integer $n \geq 1$ let $\Omega = \mathbb{R}^n$. Any sign vector $\xi \in \{-1, 1\}^n$ induces a measurable bijection $g_\xi : \mathbb{R}^n \rightarrow \mathbb{R}^n$,

$$\omega = (\omega_i)_{i=1}^n \mapsto g_\xi(\omega) := (\xi_i \omega_i)_{i=1}^n.$$

Note that $\{-1, 1\}^n$ with component-wise multiplication is an Abelian group. The corresponding family $\mathcal{G} := \{g_\xi, \xi \in \{-1, 1\}^n\}$ is an Abelian group, too. Indeed, for arbitrary $\xi, \zeta \in \{-1, 1\}^n$,

$$g_\xi \circ g_\zeta = g_\gamma \quad \text{with} \quad \gamma = (\xi_i \zeta_i)_{i=1}^n.$$

A distribution \mathbb{P} on \mathbb{R}^n is called *sign-invariant*, if it is \mathcal{G} -invariant. In other words, if $X = (X_i)_{i=1}^n$ has distribution \mathbb{P} , then for any sign vector $\xi \in \{-1, 1\}^n$, the random vector $(\xi_i X_i)_{i=1}^n$ has distribution \mathbb{P} , too.

If $\mathcal{E} = (\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (\mathbb{P}_\theta)_{\theta \in \Theta})$ consists of sign-invariant distributions, then the sub- σ -field $\mathcal{A}_{\mathcal{G}}$ of all sign-invariant sets is sufficient for \mathcal{E} . The corresponding stochastic kernel K from $(\Omega, \mathcal{A}_{\mathcal{G}})$ to (Ω, \mathcal{A}) is given by

$$K(\omega, \cdot) = \frac{1}{2^n} \sum_{\xi \in \{-1, 1\}^n} \delta_{g_\xi(\omega)}.$$

8.4 Complete Statistical Experiments

Definition 8.28 (Complete statistical experiment). A statistical experiment $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ is called *boundedly complete*, if for any bounded and measurable function $f : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}$,

$$\int f d\mathbb{P}_\theta = 0 \quad \text{for all } \theta \in \Theta$$

implies that

$$\mathbb{P}_\theta(f \neq 0) = 0 \quad \text{for all } \theta \in \Theta.$$

The experiment \mathcal{E} is called *complete* if for any function $f \in \bigcap_{\theta \in \Theta} \mathcal{L}^1(\mathbb{P}_\theta)$,

$$\int f d\mathbb{P}_\theta = 0 \quad \text{for all } \theta \in \Theta$$

implies that

$$\mathbb{P}_\theta(f \neq 0) = 0 \quad \text{for all } \theta \in \Theta.$$

Obviously, completeness of \mathcal{E} implies bounded completeness.

Example 8.29 (Simple d -parameter exponential families). Let M be a σ -finite measure on \mathbb{R}^d and $h : \mathbb{R}^d \rightarrow [0, \infty)$ a measurable function such that for all parameters θ in a set $\Theta \subset \mathbb{R}^d$,

$$0 < \int \exp(\theta^\top x) h(x) M(dx) < \infty.$$

Then we consider the statistical experiment $\mathcal{E} = (\mathbb{R}^d, \text{Borel}(\mathbb{R}^d), (\mathbb{P}_\theta)_{\theta \in \mathbb{R}^d})$ where

$$\frac{d\mathbb{P}_\theta}{dM}(x) = C(\theta) \exp(\theta^\top x) h(x)$$

with $C(\theta) := (\int \exp(\theta^\top x) h(x) M(dx))^{-1}$. If Θ has nonempty interior, then \mathcal{E} is complete.

This follows immediately from Theorem A.15 in Appendix A.3.

Example 8.30 (Product measures). Let $(\mathcal{X}, \mathcal{B}, M)$ be a σ -finite measure space, and let Θ be a family of probability densities with respect to M . For some integer $n \geq 2$ consider $\mathcal{E} = (\mathcal{X}^n, \mathcal{B}^{\otimes n}, (Q_\theta^{\otimes n})_{\theta \in \Theta})$ with Q_θ given by $dQ_\theta/dM = \theta$. This experiment is *not* boundedly complete. To see this, consider some bounded measurable function $h : \mathcal{X} \rightarrow \mathbb{R}$ with $\text{Var}_{\theta_o}(h) > 0$ for at least one $\theta_o \in \Theta$. Then $f(\omega) := h(\omega_1) - h(\omega_2)$ defines a bounded and measurable function $f : \mathcal{X}^n \rightarrow \mathbb{R}$ such that $\int f dQ_\theta^{\otimes n} = 0$ for any $\theta \in \Theta$, but $Q_{\theta_o}^{\otimes n}(f \neq 0) > 0$, because $\int f^2 dQ_{\theta_o}^{\otimes n} = 2 \text{Var}_{\theta_o}(h) > 0$.

However, if we replace $\mathcal{A} = \mathcal{B}^{\otimes n}$ with the sub- σ -field \mathcal{A}_o of all permutation invariant sets $A \in \mathcal{A}$, then the experiment $\mathcal{E}_o = (\mathcal{X}^n, \mathcal{A}_o, (Q_\theta^{\otimes n})_{\theta \in \Theta})$ is often boundedly complete or even complete:

Special case 1. Suppose that Θ contains the convex hull of all probability densities $M(B)^{-1}1_B$ with $B \in \mathcal{B}$ satisfying $0 < M(B) < \infty$. Then \mathcal{E}_o is boundedly complete.

Proof: For arbitrary fixed sets $B_1, \dots, B_n \in \mathcal{B}$ with $0 < M(B_i) < \infty$ and arbitrary tuples $\gamma \in (0, \infty)^n$, consider the probability density

$$\theta_\gamma := C_\gamma^{-1} \sum_{i=1}^n \gamma_i 1_{B_i}$$

with $C_\gamma := \sum_{i=1}^n \gamma_i M(B_i)$. Then by assumption, $\theta_\gamma \in \Theta$. If $f : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}$ is bounded and measurable, it follows from

$$\int f dQ_\theta^{\otimes n} = 0 \quad \text{for all } \theta \in \Theta$$

that

$$(8.3) \quad \int f(x) \prod_{i=1}^n \left(\sum_{j=1}^n \gamma_j 1_{B_j}(x_i) \right) M^{\otimes n}(dx) = 0 \quad \text{for all } \gamma \in (0, \infty)^n.$$

The integral on the left hand side equals

$$\sum_{j(1), \dots, j(n)=1}^n \prod_{i=1}^n \gamma_{j(i)} \int_{B_{j(1)} \times \dots \times B_{j(n)}} f dM^{\otimes n},$$

and this is an n -variate polynomial in γ of degree n . Since

$$\frac{\partial^n}{\partial \gamma_1 \partial \gamma_2 \cdots \partial \gamma_n} \prod_{i=1}^n \gamma_{j(i)} = \begin{cases} 1 & \text{if } \{j(1), \dots, j(n)\} = \{1, \dots, n\} \\ 0 & \text{else,} \end{cases}$$

it follows from (8.3) that

$$\sum_{\pi \in \mathcal{S}_n} \int_{B_{\pi(1)} \times \cdots \times B_{\pi(n)}} f dM^{\otimes n} = 0.$$

If f is permutation-invariant, that means, \mathcal{A}_o -measurable rather than just \mathcal{A} -measurable, then the latter sum equals

$$n! \int_{B_1 \times \cdots \times B_n} f dM^{\otimes n}.$$

All in all, we know that

$$\int_{B_1 \times \cdots \times B_n} f dM^{\otimes n} = 0$$

for arbitrary sets $B_1, \dots, B_n \in \mathcal{B}$ with finite measure $M(B_i)$. As shown in Exercise 8.31, this implies that $M^{\otimes n}(f \neq 0) = 0$. \square

Exercise 8.31. Let M be a σ -finite measure on $(\mathcal{X}, \mathcal{B})$, and let $f : \mathcal{X}^n \rightarrow \mathbb{R}$ be bounded and measurable such that

$$\int_{B_1 \times \cdots \times B_n} f dM^{\otimes n} = 0$$

for arbitrary sets $B_1, \dots, B_n \in \mathcal{B}$ with finite measure $M(B_i)$. Show that

$$M^{\otimes n}(f \neq 0) = 0.$$

Hint: Consider first the case of $M(\mathcal{X}) < \infty$ and the two finite measures Q^\pm given by $Q^\pm(A) := \int_A f^\pm dM^{\otimes n}$.

Special case 2. Let M be Lebesgue measure on $\mathcal{X} = \mathbb{R}$. Suppose that for an arbitrary fixed $\sigma > 0$, the set Θ contains all finite convex combinations of Gaussian densities $\phi_{\mu, \sigma}$, $\mu \in \mathbb{R}$, where

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

Then \mathcal{E}_o is complete.

Proof: With the same arguments as in the previous special case, one can show that for any \mathcal{A}_o -measurable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, it follows from

$$\int f dQ_\theta^{\otimes n} = 0 \quad \text{for arbitrary } \theta \in \Theta$$

that

$$\int f(\mathbf{x}) \prod_{i=1}^n \phi_{\mu_i, \sigma}(x_i) d\mathbf{x} = 0 \quad \text{for arbitrary } \boldsymbol{\mu} \in \mathbb{R}^n,$$

and this is equivalent to

$$\int f(\mathbf{x}) \exp(-\|\mathbf{x}\|^2/(2\sigma^2)) \exp(\boldsymbol{\mu}^\top \mathbf{x}) d\mathbf{x} = 0 \quad \text{for arbitrary } \boldsymbol{\mu} \in \mathbb{R}^n.$$

But then it follows from Theorem A.15 in the Appendix that the two finite measures $A \mapsto \int_A f^\pm dM^{\otimes n}$ are identical, whence $M^{\otimes n}(f \neq 0) = 0$. \square

Unbiased estimation. With the concepts of sufficiency and completeness one can say something about unbiased point estimators. This was first noted by P.R. Halmos [7].

Theorem 8.32 (Halmos). *Let $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ be a statistical experiment with a sufficient statistic $T : (\Omega, \mathcal{A}) \rightarrow (\tilde{\Omega}, \tilde{\mathcal{A}})$. Consider some function $g : \Theta \rightarrow \mathbb{R}^q$ and a loss function $L : \mathbb{R}^q \times \Theta \rightarrow [0, \infty)$ such that $L(\cdot, \theta)$ is convex for any fixed $\theta \in \Theta$.*

(a) *Suppose $\hat{g} : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}^q$ is an unbiased estimator of $g(\theta)$. That means,*

$$\mathbb{E}_\theta(\hat{g}) = g(\theta) \quad \text{for any } \theta \in \Theta.$$

Then

$$\check{g}(\omega) := \mathbb{E}(\hat{g} | T = T(\omega))$$

defines another unbiased estimator with

$$\mathbb{E}_\theta L(\check{g}, \theta) \leq \mathbb{E}_\theta L(\hat{g}, \theta) \quad \text{for any } \theta \in \Theta.$$

If $L(\cdot, \theta)$ is strictly convex, then the latter inequality is strict unless $\mathbb{P}_\theta(\hat{g} \neq \check{g}) = 0$.

(b) *If $\mathcal{E}^T = (\tilde{\Omega}, \tilde{\mathcal{A}}, (\mathbb{P}_\theta^T)_{\theta \in \Theta})$ is complete, then the unbiased estimator \check{g} in part (a) is essentially unique. That means, if $\tilde{g} = \tilde{h}(T)$ is another unbiased estimator of $g(\theta)$, then $\mathbb{P}_\theta(\tilde{g} \neq \check{g}) = 0$ for arbitrary $\theta \in \Theta$.*

Proof of Theorem 8.32. By sufficiency of T , there exists a stochastic kernel K from $(\tilde{\Omega}, \tilde{\mathcal{A}})$ to (Ω, \mathcal{A}) describing the conditional distribution of $X \sim \mathbb{P}_\theta$, given $T(X)$, for any parameter $\theta \in \Theta$. This gives rise to conditional expectations $\mathbb{E}(h | T = t) := \int h(\omega) K(t, d\omega)$. Note that by assumption, for any $\theta \in \Theta$,

$$\infty > \mathbb{E}_\theta \|\hat{g}\| = \int_{\tilde{\Omega}} \mathbb{E}(\|\hat{g}\| | T = t) \mathbb{P}_\theta^T(dt),$$

so the set $\tilde{N} := \{t \in \tilde{\Omega} : \mathbb{E}(\|\hat{g}\| | T = t) = \infty\}$ satisfies $\mathbb{P}_\theta^T(\tilde{N}) = 0$. In particular,

$$\check{h}(t) := \mathbb{E}(\hat{g} | T = t)$$

is well-defined in \mathbb{R}^d for any $t \in \tilde{\Omega} \setminus \tilde{N}$.

As to part (a), $\check{g} = \check{h}(T)$ with $\check{h}(t)$ defined before. Since \hat{g} is unbiased,

$$g(\theta) = \mathbb{E}_\theta(\hat{g}) = \int_{\tilde{\Omega}} \mathbb{E}(\hat{g} | T = t) \mathbb{P}_\theta^T(dt) = \int_{\tilde{\Omega}} \check{h}(t) \mathbb{P}_\theta^T(dt) = \mathbb{E}_\theta(\check{h}(T)) = \mathbb{E}_\theta(\check{g}).$$

Hence, \check{g} is unbiased, too. Moreover, applying Jensen's inequality to the conditional expectations $\mathbb{E}(\cdot | T = t)$ leads to

$$\begin{aligned} \mathbb{E}_\theta L(\hat{g}, \theta) &= \int_{\tilde{\Omega}} \mathbb{E}(L(\hat{g}, \theta) | T = t) \mathbb{P}_\theta^T(dt) \\ &\geq \int_{\tilde{\Omega}} L(\mathbb{E}(\hat{g} | T = t), \theta) \mathbb{P}_\theta^T(dt) \\ &= \int_{\tilde{\Omega}} L(\check{h}(t), \theta) \mathbb{P}_\theta^T(dt) \\ &= \mathbb{E}_\theta L(\check{g}, \theta). \end{aligned}$$

If $L(\cdot, \theta)$ is strictly convex, the inequality

$$\mathbb{E}(L(\hat{g}, \theta) | T = t) \geq L(\mathbb{E}(\hat{g} | T = t), \theta)$$

is strict, unless $P(\hat{g} \neq \check{g} | T = t) = 0$. Hence $\mathbb{E}_\theta L(\hat{g}, \theta) > \mathbb{E}_\theta L(\check{g}, \theta)$, unless $\mathbb{P}_\theta(\hat{g} \neq \check{g}) = 0$.

As to part (b), suppose that $\tilde{g} = \tilde{h}(T)$ is another unbiased estimator of $g(\theta)$. That means, the difference $\Delta := \tilde{h} - \check{h}$ satisfies

$$\int_{\tilde{\Omega}} \Delta d\mathbb{P}_\theta^T = 0 \quad \text{for all } \theta \in \Theta.$$

But then completeness of \mathcal{E}^T implies that $0 = \mathbb{P}_\theta^T(\Delta \neq 0) = \mathbb{P}_\theta(\tilde{g} \neq \check{g})$ for arbitrary $\theta \in \Theta$. \square

8.5 U-Statistics

The material in this section is based on the famous paper [8] by W. Hoeffding. Let X_1, X_2, \dots, X_n be independent random variables with unknown distribution P on a measurable space $(\mathcal{X}, \mathcal{B})$. To establish a link to the previous sections, let \mathcal{P} be a given family of probability distributions on $(\mathcal{X}, \mathcal{B})$. Assuming that the unknown distribution P belongs to \mathcal{P} , the corresponding statistical experiment is

$$(\mathcal{X}^n, \mathcal{B}^{\otimes n}, (P^{\otimes n})_{P \in \mathcal{P}}).$$

Note that all distributions $P^{\otimes n}$, $P \in \mathcal{P}$, are exchangeable (permutation-invariant). Hence for any given function $g : \mathcal{P} \rightarrow \mathbb{R}^q$, an unbiased point estimator \hat{g} of $g(P)$ can be improved by replacing $\hat{g}(X_1, X_2, \dots, X_n)$ with

$$\frac{1}{n!} \sum_{\pi \in \mathcal{S}_n} \hat{g}(X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(n)}).$$

And if \mathcal{P} is sufficiently rich, the latter estimator is essentially unique, see Example 8.30 and Theorem 8.32.

U-Statistics. Now we consider a particular type of parameter $g(P)$. Let $h : (\mathcal{X}^m, \mathcal{B}^{\otimes m}) \rightarrow \mathbb{R}$ be a given measurable function such that

$$g(P) := \int h dP^{\otimes m} \in \mathbb{R}.$$

In case of $m \geq 2$, we assume without loss of generality that h is permutation-invariant, i.e. symmetric in its m arguments. Otherwise we could replace $h(x_1, \dots, x_m)$ with

$$\tilde{h}(x_1, \dots, x_m) := \frac{1}{m!} \sum_{\pi \in \mathcal{S}_m} h(x_{\pi(1)}, \dots, x_{\pi(m)})$$

which defines a permutation-invariant function \tilde{h} satisfying $\int \tilde{h} dP^{\otimes m} = \int h dP^{\otimes m}$.

Obviously, for $n \geq m$, a naive unbiased estimator for $g(P)$ is given by $\hat{g} := h(X_1, \dots, X_m)$. Averaging this naive estimator over all permutations of $\mathbf{X}_n = (X_i)_{i=1}^n$ yields the following unbiased estimator of $g(P)$:

$$\check{g}_n := \binom{n}{m}^{-1} \sum_{1 \leq i_1 < i_2 < \dots < i_m \leq n} h(X_{i_1}, X_{i_2}, \dots, X_{i_m}).$$

Such an estimator is called a U -statistic with kernel h or a U -statistic of order m .

Example 8.33 (Mean and variance of a distribution and sample). Let $\mathcal{X} = \mathbb{R}$, and suppose that $\int |x| P(dx) < \infty$. Then the mean $\mu(P) := \int x P(dx)$ equals $g(P)$ with $h : \mathbb{R} \rightarrow \mathbb{R}$, $h(x) = x$, and the sample mean is a U -statistic of order 1:

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i = \binom{n}{1}^{-1} \sum_{1 \leq i \leq n} h(X_i).$$

In case of $\int x^2 P(dx) < \infty$ one may write

$$\sigma^2(P) := \int (x - \mu(P))^2 P(dx) = \int_{\mathbb{R}^2} h dP^{\otimes 2} \quad \text{with} \quad h(x_1, x_2) := (x_1 - x_2)^2/2.$$

The corresponding U -statistic is just the usual sample variance:

$$\begin{aligned} \binom{n}{2}^{-1} \sum_{1 \leq i_1 < i_2 \leq n} (X_{i_1} - X_{i_2})^2/2 &= \frac{1}{n(n-1)} \sum_{i,j=1}^n (X_i - X_j)^2/2 \\ &= \frac{1}{n(n-1)} \left(n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2 \right) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n \bar{X}_n^2 \right) \\ &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2. \end{aligned}$$

Example 8.34. Let $\mathcal{X} = [0, \infty)$, and suppose that we want to quantify whether P has strong right tails. On possibility to quantify this property would be consider

$$g(P) := \mathbb{P}\left(X_3 > \sqrt{X_1^2 + X_2^2}\right).$$

(Maybe this is not such a brilliant proposal; the main point is to illustrate the construction of U -statistics.) This corresponds to $\mathbb{E} h(X_1, X_2, X_3)$ with $h(x_1, x_2, x_3) := 1_{[x_3 > \sqrt{x_1^2 + x_2^2}]}$. Symmetrizing this kernel leads to

$$h(x_1, x_2, x_3) = \frac{1}{3} \sum_{i=1}^3 1_{[x_i > \sqrt{\|x\|^2 - x_i^2}]},$$

where $\|x\|^2 = x_1^2 + x_2^2 + x_3^2$, and the resulting unbiased estimator for $g(P)$ would be the U -statistic with this symmetric kernel h .

Here is a first result about the variance of a U -statistic which will suffice for our purposes:

Lemma 8.35 (Hoeffding). Suppose that $\int h^2 dP^{\otimes m} < \infty$. Let

$$\begin{aligned} h_m(x_1, \dots, x_m) &:= h(x_1, \dots, x_m) - g(P), \\ h_k(x_1, \dots, x_k) &:= \mathbb{E} h_m(x_1, \dots, x_k, X_{k+1}, \dots, X_m) \quad \text{for } 1 \leq k < m, \\ \sigma_k^2 &:= \mathbb{E}(h_k(X_1, \dots, X_k)^2) \quad \text{for } 1 \leq k \leq m. \end{aligned}$$

Then $\sigma_1^2 \leq \dots \leq \sigma_m^2$, and

$$\text{Var}(\check{g}_n) = \sum_{k=1}^m \mathbb{P}(Y = k) \sigma_k^2$$

with $Y \sim \text{Hyp}(n, m, m)$, i.e. $\mathbb{P}(Y = k) = \binom{m}{k} \binom{n-m}{m-k} / \binom{n}{m}$.

This lemma is only a simplified version of Hoeffding's [8] findings. Stronger statements about the variances σ_k^2 and the distribution of \check{g}_n will be derived in Exercise 8.43 from the so-called Hoeffding decomposition, introduced in Section A.4.

Proof of Lemma 8.35. For any index set $J \subset \{1, \dots, n\}$ with $1 \leq \#J \leq m$ let $X_J := (X_i)_{i \in J}$. Then we may write

$$\check{g}_n - g(P) = U_n := \binom{n}{m}^{-1} \sum_{J \subset \{1, \dots, n\}: \#J=m} h_m(X_J).$$

It follows from independence of the X_i , Fubini's theorem, symmetry of h and the definition of h_k , $1 \leq k \leq m$, that

$$(8.4) \quad \mathbb{E}(h_m(X_I)h_m(X_J)) = \sigma_{\#(I \cap J)}^2 \quad \text{for } I, J \subset \{1, \dots, n\}, \#I = \#J = m,$$

where $\sigma_0^2 := 0$. Indeed, $\mathbb{E} h_m(X_1, \dots, X_m) = 0$ by definition of h_m , so in case of $I \cap J = \emptyset$, it follows from stochastic independence of X_I and X_J that

$$\mathbb{E}(h_m(X_I)h_m(X_J)) = \mathbb{E}(h_m(X_I)) \mathbb{E}(h_m(X_J)) = 0.$$

Moreover, by definition of σ_m^2 ,

$$\mathbb{E}(h_m(X_I)h_m(X_J)) = \sigma_m^2 \quad \text{if } I = J.$$

If $1 \leq k := \#(I \cap J) < m$, then stochastic independence of $X_{I \cap J}$, $X_{I \setminus J}$, $X_{J \setminus I}$ and the definition of h_k imply that

$$\begin{aligned} \mathbb{E}(h_m(X_I)h_m(X_J)) &= \mathbb{E}(h_m(X_{I \cap J}, X_{I \setminus J})h_m(X_{I \cap J}, X_{J \setminus I})) \\ &= \mathbb{E} \mathbb{E}(h_m(X_{I \cap J}, X_{I \setminus J})h_m(X_{I \cap J}, X_{J \setminus I}) \mid X_{I \cap J}) \\ &= \mathbb{E}(h_k(X_{I \cap J})^2) \\ &= \sigma_k^2. \end{aligned}$$

Equation (8.4) yields the specific formula for $\text{Var}(\check{g}_n)$, because

$$\text{Var}(\check{g}_n) = \mathbb{E}(U_n U_n) = \binom{n}{m}^{-1} \sum_{I \subset \{1, \dots, n\}: \#I=m} \mathbb{E}(h_m(X_I)U_n),$$

and for any fixed index set $I \subset \{1, \dots, n\}$ with m elements,

$$\begin{aligned}
\mathbb{E}(h_m(X_I)U_n) &= \binom{n}{m}^{-1} \sum_{J \subset \{1, \dots, n\}: \#J=m} \mathbb{E}(h_m(X_I)h_m(X_J)) \\
&= \binom{n}{m}^{-1} \sum_{k=1}^m \#\{J \subset \{1, \dots, n\} : \#J = m, \#(J \cap I) = k\} \sigma_k^2 \\
&= \sum_{k=1}^m \binom{n}{m}^{-1} \binom{m}{k} \binom{n-m}{m-k} \sigma_k^2 \\
&= \sum_{k=1}^m \mathbb{P}(Y = k) \sigma_k^2.
\end{aligned}$$

Since this value does not depend on the particular choice of I , this proves the asserted equation for $\text{Var}(\check{g}_n)$.

It remains to prove the inequality $\sigma_k^2 \leq \sigma_{k+1}^2$ for $1 \leq k < m$. Note first that

$$h_k(x_1, \dots, x_k) = \mathbb{E} h_{k+1}(x_1, \dots, x_k, X_{k+1}).$$

In case of $k = m - 1$ this is just the definition of $h_k(x_1, \dots, x_k)$, and otherwise it is a consequence of Fubini's theorem, because

$$\begin{aligned}
h_k(x_1, \dots, x_k) &= \mathbb{E} h_m(x_1, \dots, x_k, X_{k+1}, X_{k+2}, \dots, X_m) \\
&= \mathbb{E} \mathbb{E}(h_m(x_1, \dots, x_k, X_{k+1}, X_{k+2}, \dots, X_m) \mid X_{k+1}) \\
&= \mathbb{E} h_{k+1}(x_1, \dots, x_k, X_{k+1}).
\end{aligned}$$

But then the Cauchy–Schwarz inequality implies that

$$\begin{aligned}
\sigma_k^2 &= \mathbb{E}(\mathbb{E}(h_{k+1}(X_1, \dots, X_k, X_{k+1}) \mid X_1, \dots, X_k)^2) \\
&\leq \mathbb{E}(\mathbb{E}(h_{k+1}(X_1, \dots, X_k, X_{k+1})^2 \mid X_1, \dots, X_k)) \\
&= \mathbb{E}(h_{k+1}(X_1, \dots, X_k, X_{k+1})^2) \\
&= \sigma_{k+1}^2.
\end{aligned}$$

□

Corollary 8.36. *In case of $\int h^2 dP^{\otimes m} < \infty$,*

$$\text{Var}(\check{g}_n) = \frac{m^2 \sigma_1^2}{n} + O(n^{-2}).$$

Proof of Corollary 8.36. This expansion follows immediately from Lemma 8.35 and the fact that

for $Y \sim \text{Hyp}(n, m, m)$,

$$\begin{aligned}
\mathbb{P}(Y = k) &= \binom{n}{m}^{-1} \binom{m}{k} \binom{n-m}{m-k} \\
&= \frac{m! [m]_k [n-m]_{m-k}}{[n]_m k! (m-k)!} \\
&= \frac{[m]_k^2 n^{m-k} (1 + O(n^{-1}))}{k! n^m (1 + O(n^{-1}))} \\
&= \frac{[m]_k^2 (1 + O(n^{-1}))}{k! n^k} \\
&= \begin{cases} m^2/n + O(n^{-2}) & \text{if } k = 1, \\ O(n^{-2}) & \text{if } k \geq 2. \end{cases}
\end{aligned}$$

□

The following result shows that U -statistics may be approximated by an average of independent, identically distributed random variables and satisfy a Central Limit Theorem:

Theorem 8.37 (Hoeffding). *Under the conditions of Lemma 8.35,*

$$\check{g}_n = g(P) + \frac{m}{n} \sum_{i=1}^n h_1(X_i) + R_n$$

where

$$\mathbb{E}(R_n^2) = O(n^{-2}).$$

Moreover,

$$\sqrt{n}(\check{g}_n - g(P)) \rightarrow_{\mathcal{L}} \mathcal{N}(0, m^2 \sigma_1^2)$$

as $n \rightarrow \infty$.

Our proof of Theorem 8.37 utilizes a general approximation result of Jaroslav Hájek:

Lemma 8.38 (Hájek projection). *Let X_1, X_2, \dots, X_n be arbitrary independent random variables with values in $(\mathcal{X}_1, \mathcal{B}_1), (\mathcal{X}_2, \mathcal{B}_2), \dots, (\mathcal{X}_n, \mathcal{B}_n)$, and let T be a random variable of the form $T = f(X_1, X_2, \dots, X_n)$ such that $\mathbb{E}T = 0$ and $\mathbb{E}(T^2) < \infty$. Then $T_i := \mathbb{E}(T | X_i)$ satisfies $\mathbb{E}(T_i) = 0$ for $1 \leq i \leq n$. Moreover, for arbitrary random variables Y_1, Y_2, \dots, Y_n of type $Y_i = f_i(X_i)$ with $\mathbb{E}(Y_i) = 0$ and $\mathbb{E}(Y_i^2) < \infty$,*

$$\begin{aligned}
\mathbb{E}\left(\left(T - \sum_{i=1}^n Y_i\right)^2\right) &= \mathbb{E}(T^2) - \sum_{i=1}^n \mathbb{E}(T_i^2) + \sum_{i=1}^n \mathbb{E}((Y_i - T_i)^2) \\
&\geq \mathbb{E}(T^2) - \sum_{i=1}^n \mathbb{E}(T_i^2)
\end{aligned}$$

with equality if and only if $Y_i = T_i$ almost surely for $1 \leq i \leq n$.

Proof of Lemma 8.38. By Fubini's theorem, $0 = \mathbb{E}T = \mathbb{E} \mathbb{E}(T | X_i) = \mathbb{E}T_i$. Moreover,

$$\begin{aligned}
\mathbb{E}\left(\left(T - \sum_{i=1}^n Y_i\right)^2\right) &= \mathbb{E}(T^2) - 2 \sum_{i=1}^n \mathbb{E}(TY_i) + \sum_{i,j=1}^n \mathbb{E}(Y_i Y_j) \\
&= \mathbb{E}(T^2) - 2 \sum_{i=1}^n \mathbb{E} \mathbb{E}(TY_i | X_i) + \sum_{i=1}^n \mathbb{E}(Y_i^2) \\
&= \mathbb{E}(T^2) - 2 \sum_{i=1}^n \mathbb{E}(\mathbb{E}(T | X_i) Y_i) + \sum_{i=1}^n \mathbb{E}(Y_i^2) \\
&= \mathbb{E}(T^2) - 2 \sum_{i=1}^n \mathbb{E}(T_i Y_i) + \sum_{i=1}^n \mathbb{E}(Y_i^2) \\
&= \mathbb{E}(T^2) - \sum_{i=1}^n \mathbb{E}(T_i^2) + \sum_{i=1}^n \mathbb{E}((Y_i - T_i)^2),
\end{aligned}$$

because $\mathbb{E}(Y_i Y_j) = \mathbb{E}(Y_i) \mathbb{E}(Y_j) = 0$ in case of $i \neq j$. □

Proof of Theorem 8.37. It follows from Lemma 8.38 that

$$\check{g}_n - g(P) = U_n = \sum_{i=1}^n \mathbb{E}(U_n | X_i) + R_n$$

with

$$\mathbb{E}(R_n^2) = \mathbb{E}(U_n^2) - \sum_{i=1}^n \mathbb{E}(\mathbb{E}(U_n | X_i)^2).$$

But

$$\begin{aligned}
\mathbb{E}(U_n | X_i) &= \binom{n}{m}^{-1} \sum_{J \subset \{1, \dots, m\}: \#J=m} \mathbb{E}(h_m(X_J) | X_i) \\
&= \binom{n}{m}^{-1} \sum_{J \subset \{1, \dots, m\}: \#J=m} 1_{[i \in J]} h_1(X_i) \\
&= \binom{n}{m}^{-1} \#\{J \subset \{1, \dots, m\} : \#J = m, i \in J\} h_1(X_i) \\
&= \binom{n}{m}^{-1} \binom{n-1}{m-1} h_1(X_i) \\
&= \frac{m}{n} h_1(X_i).
\end{aligned}$$

Hence

$$U_n = \frac{m}{n} \sum_{i=1}^n h_1(X_i) + R_n$$

with

$$\mathbb{E}(R_n^2) = \mathbb{E}(U_n^2) - \frac{m^2 \sigma_1^2}{n} = \text{Var}(\check{g}_n) - \frac{m^2 \sigma_1^2}{n} = O(n^{-2})$$

according to Corollary 8.36. Consequently, $R_n = O_p(n^{-1})$, whence

$$\sqrt{n}(\check{g}_n - g(P)) = \frac{m}{\sqrt{n}} \sum_{i=1}^n h_1(X_i) + O_p(n^{-1/2}),$$

and it follows from the usual Central Limit Theorem that the right hand side converges in distribution to a Gaussian random variable with mean 0 and variance $m^2 \mathbb{E}(h_1(X_1)^2) = m^2 \sigma_1^2$. \square

Remark 8.39 (Hoeffding's decomposition in case of $m = 2$). In case of $m = 2$ one may write

$$\check{g}_n = g(P) + \frac{2}{n} \sum_{i=1}^n h_1(X_i) + \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} h_2^o(X_i, X_j)$$

with

$$\begin{aligned} h_2^o(x_1, x_2) &:= h_2(x_1, x_2) - h_1(x_1) - h_1(x_2) \\ &= h(x_1, x_2) - \mathbb{E} h(x_1, X_2) - \mathbb{E} h(X_1, x_2) + \mathbb{E} h(X_1, X_2). \end{aligned}$$

Moreover, the $n + \binom{n}{2}$ random variables $h_1(X_i)$ ($1 \leq i \leq n$) and $h_2^o(X_i, X_j)$ ($1 \leq i < j \leq n$) are easily shown to be centered and uncorrelated with

$$\mathbb{E}(h_2^o(X_i, X_j)^2) = \sigma_2^2 - 2\sigma_1^2 \leq \sigma_2^2.$$

Hence the remainder

$$R_n = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} h_2^o(X_i, X_j)$$

satisfies the (in)equalities

$$\mathbb{E}(R_n^2) = \binom{n}{2}^{-1} (\sigma_2^2 - 2\sigma_1^2) = \frac{2\sigma_2^2 - 4\sigma_1^2}{n(n-1)} \leq \frac{2\sigma_2^2}{n(n-1)}.$$

Example 8.33 (Sample variance, continued). With $h(x_1, x_2) = (x_1 - x_2)^2/2$, the auxiliary function h_1 is given by

$$\begin{aligned} h_1(x) &= \mathbb{E} h(x, X_1) - \sigma^2(P) \\ &= x^2/2 - x\mu(P) + \mathbb{E}(X_1^2)/2 - \sigma^2(P) \\ &= [(x - \mu(P))^2 - \sigma^2(P)]/2, \end{aligned}$$

and this leads to the representation

$$\begin{aligned} S_X^2 &= \binom{n}{2} \sum_{1 \leq i < j \leq n} h(X_i, X_j) \\ &= \sigma^2(P) + \frac{2}{n} \sum_{i=1}^n h_1(X_i) + R_n \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu(P))^2 + R_n, \end{aligned}$$

where

$$\mathbb{E}(R_n^2) = O(n^{-2}) \quad \text{if} \quad \int x^4 P(dx) < \infty.$$

Moreover, as $n \rightarrow \infty$,

$$\sqrt{n}(S_X^2 - \sigma(P)^2) \rightarrow_{\mathcal{L}} \mathcal{N}(0, \mathbb{E}[(X_1 - \mu(P))^4] - \sigma(P)^4).$$

Example 8.40 (Kendall's τ). Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be independent random pairs with distribution P on $\mathbb{R} \times \mathbb{R}$. A nonparametric measure of correlation of X_1 and Y_1 , proposed by Maurice Kendall [9], is given by

$$\tau(P) := \mathbb{E}(\text{sign}(X_2 - X_1) \text{sign}(Y_2 - Y_1)).$$

This is the probability, that the two observation pairs (X_1, Y_1) and (X_2, Y_2) are “concordant”, i.e.

$$\text{sign}(X_2 - X_1) = \text{sign}(Y_2 - Y_1) \neq 0,$$

minus the probability that they are “discordant”, i.e.

$$-\text{sign}(X_2 - X_1) = \text{sign}(Y_2 - Y_1) \neq 0.$$

If X_1 and Y_1 are stochastically independent, then the four random variables X_1, X_2, Y_1, Y_2 are stochastically independent, and

$$\tau(P) = \mathbb{E} \text{sign}(X_2 - X_1) \mathbb{E} \text{sign}(Y_2 - Y_1) = 0,$$

because the distributions of $X_2 - X_1$ and $Y_2 - Y_1$ are symmetric around 0.

Note that $\tau(P) = \mathbb{E} h((X_1, Y_1), (X_2, Y_2))$ with the kernel

$$h((x_1, y_1), (x_2, y_2)) := \text{sign}(x_2 - x_1) \text{sign}(y_2 - y_1).$$

Consequently, an unbiased estimator for $\tau(P)$ is given by Kendall's τ statistic

$$\check{\tau}_n := \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} \text{sign}(X_j - X_i) \text{sign}(Y_j - Y_i),$$

which is a U -statistic of order 2 with kernel h .

Since $|h| \leq 1$, we may apply Theorem 8.37 and conclude that

$$\check{\tau}_n = \tau(P) + \frac{2}{n} \sum_{i=1}^n h_1(X_i, Y_i) + R_n$$

with $h_1(x, y) := \mathbb{E}(\text{sign}(x - X_1) \text{sign}(y - Y_1)) - \tau(P)$ and a remainder term R_n such that $\mathbb{E}(R_n^2) = O(n^{-2})$. Moreover, as $n \rightarrow \infty$,

$$\sqrt{n}(\check{\tau}_n - \tau(P)) \rightarrow_{\mathcal{L}} \mathcal{N}(0, 4 \mathbb{E}(h_1(X_1, Y_1)^2)).$$

Let us consider the following special case: Suppose that X_1 and Y_1 are stochastically independent with continuous distribution function F and G , respectively. Then one can easily verify that

$$h_1(x, y) = (2F(x) - 1)(2G(y) - 1).$$

But $2F(X_1) - 1$ and $2G(Y_1) - 1$ are stochastically independent and uniformly distributed on $[-1, 1]$. From this one can easily deduce that

$$\mathbb{E}(h_1(X_1, Y_1)^2) = 1/9,$$

so

$$\sqrt{n} \check{\tau}_n \rightarrow_{\mathcal{L}} \mathcal{N}(0, 4/9).$$

Exercise 8.41. Let X_1, X_2, \dots, X_n be independent random variables with unknown distribution P on \mathbb{R} . For $m \in \mathbb{N}$ let

$$g_m(P) := \mathbb{E}_P \text{Med}(X_1, \dots, X_m)$$

with sample median function $\text{Med}(\dots)$.

Show that for $n \geq m$, the corresponding U-statistic

$$\check{g} = \binom{n}{m}^{-1} \sum_{1 \leq i_1 < \dots < i_m \leq n} \text{Med}(X_{i_1}, \dots, X_{i_m})$$

is a L-statistic, that means,

$$\check{g} = \sum_{i=1}^n w_i X_{(i)}$$

with suitable weights $w_1, w_2, \dots, w_n \geq 0$.

Hint: Distinguish the cases of odd and even m .

Exercise 8.42. Suppose that X_1, X_2, X_3, \dots are independent and identically distributed with distribution P on \mathbb{R} . Suppose further that $\mathbb{E}(|X_1|^3) < \infty$.

(a) Determine an optimal unbiased estimator of the centered third moment,

$$g(P) := \mathbb{E}((X_1 - \mathbb{E}(X_1))^3).$$

Hint: Determine a measurable function $h : \mathbb{R}^3 \rightarrow \mathbb{R}$ such that $g(P) = \mathbb{E} h(X_1, X_2, X_3)$, and construct a corresponding U-statistic \check{g} .

(b) A naive estimator for $g(P)$ is given by

$$\hat{g}_{\text{naive}} := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^3.$$

Show that this estimator can be written as a function of the sums $S_\ell := \sum_{i=1}^n X_i^\ell$, $1 \leq \ell \leq 3$, and that the computation of \hat{g}_{naive} requires $O(n)$ steps.

(c) Show that the unbiased estimator \check{g} is also a simple function of S_1, S_2, S_3 and can be computed in $O(n)$ steps.

(d) Show that for some constant $C(P)$,

$$\mathbb{E} |\check{g} - \hat{g}_{\text{naive}}| \leq C(P)n^{-1}$$

for all $n \geq 3$. Proposal: One can write $\check{g} - \hat{g}_{\text{naive}}$ as a linear combination of the three sums $S_3, S_{21} := \sum_{i,j=1}^n 1_{[i \neq j]} X_i^2 X_j$ and $S_{111} := \sum_{i,j,k=1}^n 1_{[i,j,k \text{ different}]} X_i X_j X_k$ with coefficients depending on n .

Exercise 8.43 (Refinements via Hoeffding's decomposition). With the notation of Lemma 8.35 and its proof, let

$$h_k^o(x_1, \dots, x_k) := \sum_{\ell=1}^k (-1)^{k-\ell} \sum_{I \subset \{1, \dots, k\}; \#I=\ell} h_\ell(x_I)$$

for $1 \leq k \leq m$. The general Hoeffding decomposition, presented in Section A.4, implies that the random variables

$$h_k^o(X_J), \quad k \in \{1, \dots, m\} \text{ and } J \subset \{1, \dots, n\}, \#J = k,$$

are centered and uncorrelated, where

$$h_k(x_1, \dots, x_k) = \sum_{\ell=1}^k \sum_{J \subset \{1, \dots, k\}: \#J=\ell} h_\ell^o(x_J).$$

Let

$$\tau_k^2 := \mathbb{E}(h_k^o(X_1, \dots, X_k)^2).$$

(a) Show that $\sigma_k^2 = \mathbb{E}(h_k(X_1, \dots, X_k)^2)$ equals

$$\sigma_k^2 = \sum_{\ell=1}^k \binom{k}{\ell} \tau_\ell^2.$$

Deduce from this representation that σ_k^2/k is non-decreasing in $k \in \{1, \dots, m\}$.

(b) Show that $U_n := \check{g}_n - g(P)$ equals

$$U_n = \sum_{k=1}^m \binom{n}{k}^{-1} \sum_{J \subset \{1, \dots, n\}: \#J=k} \binom{m}{k} h_k^o(X_J),$$

and that

$$\mathbb{E}(U_n^2) = \sum_{k=1}^m \frac{[m]_k}{[n]_k} \binom{m}{k} \tau_k^2.$$

(c) Deduce from (a) and (b) that

$$\frac{m^2 \sigma_1^2}{n} \leq \mathbb{E}(U_n^2) \leq \frac{m^2 \sigma_1^2}{n} + \frac{m^2 \sigma_m^2}{n^2}.$$

Show that $n \mathbb{E}(U_n^2)$ is non-increasing in $n \geq m$ with limit $m^2 \sigma_1^2$ as $n \rightarrow \infty$.

Chapter 9

Exponential Families

9.1 Definitions and Basic Properties

Definition 9.1 (Exponential families). A statistical experiment $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ is called an *exponential family* if there exist a σ -finite measure M on (Ω, \mathcal{A}) , a measurable mapping $h : \Omega \rightarrow [0, \infty)$, a mapping $\alpha : \Theta \rightarrow \mathbb{R}^d$ and a measurable mapping $T : \Omega \rightarrow \mathbb{R}^d$ such that for any $\theta \in \Theta$,

$$\frac{d\mathbb{P}_\theta}{dM}(\omega) = h(\omega) \exp(\alpha(\theta)^\top T(\omega) - \kappa(\theta))$$

with

$$\kappa(\theta) := \log \int h \exp(\alpha(\theta)^\top T) dM.$$

(In particular, we assume that $\int h \exp(\alpha(\theta)^\top T) dM < \infty$ for all $\theta \in \Theta$.)

Definition 9.2 (Natural exponential families). For a given measure space (Ω, \mathcal{A}, M) and measurable mappings $h : \Omega \rightarrow [0, \infty)$, $T : \Omega \rightarrow \mathbb{R}^d$, the corresponding *natural exponential family* is given by $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta_{\text{nat}}})$ with the *natural parameter space*

$$\Theta_{\text{nat}} := \left\{ \theta \in \mathbb{R}^d : \int h \exp(\theta^\top T) dM < \infty \right\},$$

and the probability distributions \mathbb{P}_θ are given by

$$\begin{aligned} \frac{d\mathbb{P}_\theta}{dM}(\omega) &:= h(\omega) \exp(\theta^\top T(\omega) - \kappa(\theta)), \\ \kappa(\theta) &:= \log \int h \exp(\theta^\top T) dM. \end{aligned}$$

Example 9.3 (Gaussian samples). Let $\Omega = \mathbb{R}^n$, equipped with its Borel σ -field, let $\Theta = \mathbb{R} \times (0, \infty)$, and for $\theta = (\mu, \sigma) \in \Theta$, let

$$\mathbb{P}_{\mu, \sigma} := \mathcal{N}(\mu, \sigma^2)^{\otimes n}.$$

With M denoting Lebesgue measure on \mathbb{R}^n , $h(\omega) := (2\pi)^{-n/2}$ and $X_i(\omega) := \omega_i$ for $\omega \in \Omega$, the

log-density of $\mathbb{P}_{\mu,\sigma}$ with respect to M is given by

$$\begin{aligned} \frac{d\mathbb{P}_{\mu,\sigma}}{dM} &= h(\omega) \exp\left(-\sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2} - n \log \sigma\right) \\ &= h(\omega) \exp\left(\frac{\mu}{\sigma^2} \sum_{i=1}^n X_i + \frac{-1}{2\sigma^2} \sum_{i=1}^n X_i^2 - n \log \sigma\right) \\ &= h(\omega) \exp(\alpha(\mu, \sigma)^\top T - \kappa(\mu, \sigma)) \end{aligned}$$

with

$$\begin{aligned} \alpha(\mu, \sigma) &:= \left(\frac{\mu}{\sigma^2}, \frac{-1}{2\sigma^2}\right)^\top, \\ T &:= \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2\right)^\top, \\ \kappa(\mu, \sigma) &:= n \log \sigma. \end{aligned}$$

Here one easily verifies that

$$\Theta_{\text{nat}} = \mathbb{R} \times (-\infty, 0).$$

Example 9.4 (Gamma distributions). Let $\Omega = (0, \infty)$, equipped with its Borel σ -field. For $a, b > 0$ let $\text{Gamma}(a, b)$ be the gamma distribution on $(0, \infty)$ with shape parameter $a > 0$ and scale parameter $b > 0$. That means, $\text{Gamma}(a, b)$ has Lebesgue density

$$\begin{aligned} f_{a,b}(\omega) &= \frac{1}{\Gamma(a)b} \left(\frac{\omega}{b}\right)^{a-1} \exp\left(-\frac{\omega}{b}\right) \\ &= \exp(\alpha(a, b)^\top T(\omega) - \kappa(a, b)) \end{aligned}$$

with

$$\begin{aligned} \alpha(a, b) &:= (a - 1, -1/b)^\top, \\ T(\omega) &:= (\log \omega, \omega)^\top, \\ \kappa(a, b) &:= a \log b + \log \Gamma(a). \end{aligned}$$

Here one easily verifies that

$$\Theta_{\text{nat}} = (-1, \infty) \times (-\infty, 0).$$

Remark 9.5 (Convexity and smoothness). Let \mathcal{E} be a natural exponential family as in Definition 9.2. It follows from convexity of the exponential function that the set Θ_{nat} is a convex subset of \mathbb{R}^d . Moreover, if $f \in \bigcap_{\theta \in \Theta_{\text{nat}}} L^1(\mathbb{P}_\theta)$, then the function $h : \Theta \rightarrow \mathbb{R}$,

$$h(\theta) := \int f d\mathbb{P}_\theta,$$

is twice continuously differentiable on the interior of Θ_{nat} with gradient

$$\nabla h(\theta) = \text{Cov}_\theta(f, T) = \int f T d\mathbb{P}_\theta - h(\theta) \int T d\mathbb{P}_\theta.$$

Remark 9.6 (Sufficiency). Suppose that \mathcal{E} is an exponential family as in Definition 9.1, where (Ω, d) is a separable and complete metric space and $\mathcal{A} = \text{Borel}(\Omega, d)$. Then T is a sufficient statistic for \mathcal{E} . This follows immediately from Neyman's factorization criterion.

Theorem 9.7 (Completeness in exponential families). *Let \mathcal{E} be an exponential family as in Definition 9.1. Suppose that the set $\{\alpha(\theta) : \theta \in \Theta\} \subset \mathbb{R}^d$ contains an interior point. Then the statistical model $(\mathbb{R}^d, \text{Borel}(\mathbb{R}^d), (\mathbb{P}_\theta^T)_{\theta \in \Theta})$ is complete.*

Proof of Theorem 9.7. We may assume without loss of generality that $h \equiv 1$. Otherwise we could replace M with \tilde{M} , where $\tilde{M}(B) := \int_B h dM$. Note that the image measures¹ M^T and \mathbb{P}_θ^T on \mathbb{R}^d satisfy

$$\frac{d\mathbb{P}_\theta^T}{dM^T}(x) = \exp(\alpha(\theta)^\top x - \kappa(\theta)).$$

Hence for a measurable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the property

$$\int_{\mathbb{R}^d} f d\mathbb{P}_\theta^T = 0 \quad \text{for all } \theta \in \Theta$$

is equivalent to

$$\int_{\mathbb{R}^d} f(x) \exp(\alpha(\theta)^\top x) M^T(dx) = 0 \quad \text{for all } \theta \in \Theta.$$

Since $\alpha(\Theta)$ contains an interior point, it follows from Theorem A.15 in Section A.3 that $f(x) = 0$ for M^T -almost all $x \in \mathbb{R}^d$. In particular, $\mathbb{P}_\theta^T(f \neq 0) = 0$ for all $\theta \in \Theta$. \square

9.2 Nuisance Parameters

In this section we consider statistical experiments $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ of the following type: Θ is a convex open subset of $\mathbb{R}^d \times \mathbb{R}$ with projections

$$\begin{aligned} N &:= \{\nu \in \mathbb{R}^d : (\nu, \gamma) \in \Theta \text{ for some } \gamma \in \mathbb{R}\}, \\ \Gamma &:= \{\gamma \in \mathbb{R} : (\nu, \gamma) \in \Theta \text{ for some } \nu \in \mathbb{R}^d\}. \end{aligned}$$

Each parameter $\theta = (\nu, \gamma) \in \Theta$ consists of a “nuisance parameter” $\nu \in N$ and a parameter $\gamma \in \Gamma$ of primary interest. The question is how to deal with the nuisance parameter ν if we are only interested in γ .

We assume that \mathcal{E} is an exponential family with natural parametrization: There exist a σ -finite measure M on (Ω, \mathcal{A}) and measurable functions $S : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}^d$, $Y : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}$ such that for arbitrary $(\nu, \gamma) \in \Theta$,

$$\frac{d\mathbb{P}_{\nu, \gamma}}{dM} = \exp(\nu^\top S + \gamma Y - \kappa(\nu, \gamma))$$

with

$$\kappa(\nu, \gamma) := \log \int \exp(\nu^\top S + \gamma Y) dM < \infty.$$

Here the pair (S, Y) is a sufficient statistic, provided that (Ω, d) is a separable and complete metric space, equipped with its Borel σ -field. Thus we may restrict our attention to decision procedures depending only on (S, Y) . The following result shows that under the measure $\mathbb{P}_{\nu, \gamma}$, the conditional distribution of Y , given that $S = s$, depends on s and the parameter γ but not on the nuisance parameter ν . Hence we may get rid of the nuisance parameter by conditioning on S .

¹ $M^T(B) := M(T \in B)$ and $\mathbb{P}_\theta^T(B) := \mathbb{P}_\theta(T \in B)$

Proposition 9.8. *Let us fix any parameter $(\nu_o, \gamma_o) \in \Theta$ and choose a stochastic kernel K from \mathbb{R}^d to \mathbb{R} describing the conditional distribution of Y , given S , under the measure $\mathbb{P}_{\nu_o, \gamma_o}$. That means,*

$$\mathbb{P}_{\nu_o, \gamma_o}(Y \in B \mid S = \cdot) = K(\cdot, B) \quad \text{for any Borel set } B \subset \mathbb{R}.$$

Then there exist a measurable weight function $w : \mathbb{R} \rightarrow (0, \infty)$ and a Borel set $C \subset \mathbb{R}^d$ with the following properties: $M(S \in C) = 0$, and for any $\gamma \in \Gamma$,

$$\tilde{\kappa}(s, \gamma) := \log \int_{\mathbb{R}} \exp(\gamma y) w(y) K(s, dy) < \infty \quad \text{for all } s \in \mathbb{R}^d \setminus C.$$

Moreover,

$$K_\gamma(s, B) := \begin{cases} \int_B \exp(\gamma y - \tilde{\kappa}(s, \gamma)) w(y) K(s, dy) & \text{if } s \in \mathbb{R}^d \setminus C, \\ 1_{[0 \in B]} & \text{if } s \in C, \end{cases}$$

defines a stochastic kernel K_γ from \mathbb{R}^d to \mathbb{R} describing the conditional distribution of Y , given S , under any measure $\mathbb{P}_{\nu, \gamma}$, $\nu \in N$. That means,

$$\mathbb{P}_{\nu, \gamma}(Y \in B \mid S = \cdot) = K_\gamma(\cdot, B) \quad \text{for all } \nu \in N \text{ and any Borel set } B \subset \mathbb{R}.$$

Proof of Proposition 9.8. The assumption on K is equivalent to

$$\mathbb{E}_{\nu_o, \gamma_o}(g(S)h(Y)) = \mathbb{E}_{\nu_o, \gamma_o}\left(g(S) \int h(y) K(S, dy)\right)$$

for arbitrary measurable functions $g : \mathbb{R}^d \rightarrow [0, \infty]$ and $h : \mathbb{R} \rightarrow [0, \infty]$. For any other parameter $(\nu, \gamma) \in \Theta$, this implies that

$$(9.1) \quad \mathbb{E}_{\nu, \gamma}(g(S)h(Y)) = \mathbb{E}_{\nu_o, \gamma_o}\left(g(S)h(Y) \frac{\exp(\nu^\top S + \gamma Y - \kappa(\nu, \gamma))}{\exp(\nu_o^\top S + \gamma_o Y - \kappa(\nu_o, \gamma_o))}\right)$$

$$(9.2) \quad = \mathbb{E}_{\nu_o, \gamma_o}\left(g(S) \exp((\nu - \nu_o)^\top S - \kappa(\nu, \gamma)) \int h(y) \exp(\gamma y) w(y) K(S, dy)\right)$$

with

$$w(y) := \exp(\kappa(\nu_o, \gamma_o) - \gamma_o y).$$

Taking $h \equiv 1$ shows that

$$\mathbb{E}_{\nu, \gamma}(g(S)) = \mathbb{E}_{\nu_o, \gamma_o}(g(S) f_{\nu, \gamma}(S))$$

with

$$f_{\nu, \gamma}(s) := \exp((\nu - \nu_o)^\top s - \kappa(\nu, \gamma)) \int \exp(\gamma y) w(y) K(s, dy).$$

Taking $g \equiv 1$ shows that for any fixed $\gamma \in \Gamma$, the set $C(\gamma)$ of all $s \in \mathbb{R}^d$ such that the integral $\int \exp(\gamma y) w(y) K(s, dy)$ is infinite satisfies $\mathbb{P}_{\nu, \gamma}(S \in C(\gamma)) = 0$ for all $\nu \in N$. Since $\mathbb{P}_{\nu, \gamma}$ has a strictly positive density with respect to M , we may even conclude that

$$M(S \in C(\gamma)) = 0.$$

But then the set $C := \bigcup_{\gamma \in \Gamma \cap \mathbb{Q}} C(\gamma)$ satisfies

$$M(S \in C) = 0,$$

and by convexity of the exponential function,

$$\int \exp(\gamma y) w(y) K(s, dy) < \infty \quad \text{for all } s \in \mathbb{R}^d \setminus C \text{ and } \gamma \in \Gamma.$$

Hence the stochastic kernel K_γ described in the proposition is well-defined, and it follows from (9.2) that for arbitrary measurable functions $g : \mathbb{R}^d \rightarrow [0, \infty]$ and $h : \mathbb{R} \rightarrow [0, \infty]$,

$$\begin{aligned} \mathbb{E}_{\nu, \gamma}(g(S)h(Y)) &= \mathbb{E}_{\nu_o, \gamma_o} \left(g(S) \int h(y) K_\gamma(S, dy) f_{\nu, \gamma}(S) \right) \\ &= \mathbb{E}_{\nu, \gamma} \left(g(S) \int h(y) K_\gamma(S, dy) \right). \end{aligned}$$

This shows that the kernel K_γ describes the conditional distribution of Y , given S , under the measure $\mathbb{P}_{\nu, \gamma}$. \square

Neyman's construction of tests. Suppose we want to find a good level- α test of the null hypothesis

$$\Theta_o = \{(\nu, \gamma) : \gamma \leq \gamma_o\}$$

for a given number $\gamma_o \in \Gamma$ such that

$$\{\nu \in \mathbb{R}^d : (\nu, \gamma_o) \in \Theta\} = N.$$

To this end we fix any nuisance parameter ν_o and choose for any $s \in \mathbb{R}^d$ numbers $k_\alpha(s) \in \mathbb{R}$ and $\gamma_\alpha(s) \in [0, 1]$ such that the test $\phi_\alpha : \mathbb{R}^d \times \mathbb{R} \rightarrow [0, 1]$ with

$$\phi_\alpha(s, y) := \begin{cases} 0 & \text{if } y < k_\alpha(s) \\ \gamma_\alpha(s) & \text{if } y = k_\alpha(s) \\ 1 & \text{if } y > k_\alpha(s) \end{cases}$$

satisfies

$$\mathbb{E}_{\nu_o, \gamma_o}(\phi_\alpha(S, Y) | S) = \alpha \quad \text{almost surely.}$$

Then ϕ_α is a level- α test of Θ_o , and has a certain optimality property:

Theorem 9.9 (UMP unbiased tests). *For given test level $\alpha \in (0, 1)$, let ϕ_α be the special test just described. This test belongs to the class Φ_α of all tests $\phi : \mathbb{R}^d \times \mathbb{R} \rightarrow [0, 1]$ such that*

$$\mathbb{E}_{\nu, \gamma} \phi(S, Y) \begin{cases} \leq \alpha & \text{if } \gamma \leq \gamma_o, \\ \geq \alpha & \text{if } \gamma > \gamma_o. \end{cases}$$

For arbitrary $\phi \in \Phi_\alpha$ and $(\nu, \gamma) \in \Theta \setminus \Theta_o$,

$$\mathbb{E}_{\nu, \gamma} \phi_\alpha(S, Y) \geq \mathbb{E}_{\nu, \gamma} \phi(S, Y).$$

Proof of Theorem 9.9. With the stochastic kernels K_γ , $\gamma \in \Gamma$, and the Borel set $C \subset \mathbb{R}^d$ as in Proposition 9.8, we may alter the set C if necessary such that

$$\int_{\mathbb{R}} \phi_\alpha(s, y) K_{\gamma_o}(s, dy) = \alpha \quad \text{for all } s \in \mathbb{R}^d \setminus C.$$

For any test $\phi : \mathbb{R}^d \times \mathbb{R} \rightarrow [0, 1]$, its power

$$\mathbb{E}_{\nu, \gamma} \phi(S, Y) = \int_{\Omega} \phi(S, Y) \exp(\nu^\top S + \gamma Y - \kappa(\nu, \gamma)) dM$$

is a continuous function of $\gamma \in \Gamma(\nu)$ for any fixed $\nu \in N$; see exercises. Here $\Gamma(\nu)$ denotes the open interval

$$\Gamma(\nu) := \{\gamma \in \mathbb{R} : (\nu, \gamma) \in \Theta\} \ni \gamma_o.$$

If $\phi \in \Phi_\alpha$, this implies that

$$\mathbb{E}_{\nu, \gamma_o} \phi(S, Y) = \alpha \quad \text{for arbitrary } \nu \in N.$$

We may also write

$$\mathbb{E}_{\nu, \gamma} \phi(S, Y) = \mathbb{E}_{\nu, \gamma} \left(\int_{\mathbb{R}} \phi(S, y) K_\gamma(S, dy) \right).$$

But the restricted statistical experiment $\mathcal{E}_o := (\Omega, \mathcal{A}, (\mathbb{P}_{\nu, \gamma_o})_{\nu \in N})$ is an exponential family with natural parametrization, because

$$\frac{d\mathbb{P}_{\nu, \gamma_o}}{dM} = \exp(\nu^\top S - \kappa(\nu, \gamma_o)) \exp(\gamma_o Y),$$

i.e. with the modified measure

$$M_o(d\omega) := \exp(\gamma_o Y(\omega)) M(d\omega)$$

we may write

$$\frac{d\mathbb{P}_{\nu, \gamma_o}}{dM_o} = \exp(\nu^\top S - \kappa(\nu, \gamma_o)).$$

Since N is open, the corresponding family $(\mathbb{R}^d, \text{Borel}(\mathbb{R}^d), (\mathbb{P}_{\nu, \gamma_o}^S)_{\nu \in N})$ is complete, that means, it follows from

$$\alpha = \mathbb{E}_{\nu, \gamma_o} \left(\int_{\mathbb{R}} \phi(S, y) K_{\gamma_o}(S, dy) \right) \quad \text{for all } \nu \in N$$

that

$$\int_{\mathbb{R}} \phi(S, y) K_{\gamma_o}(S, dy) = \alpha \quad \text{almost surely}$$

under any measure $\mathbb{P}_{\nu, \gamma_o}$, $\nu \in N$. But our special construction of ϕ_α and Theorem 7.13 imply that for $\gamma > \gamma_o$ and $s \in \mathbb{R}^d \setminus C$,

$$\int_{\mathbb{R}} \phi(s, y) K_\gamma(s, dy) \leq \int_{\mathbb{R}} \phi_\alpha(s, y) K_\gamma(s, dy) \quad \text{whenever} \quad \int_{\mathbb{R}} \phi(s, y) K_{\gamma_o}(s, dy) \leq \alpha.$$

Thus for arbitrary $(\nu, \gamma) \in \Theta$ with $\gamma > \gamma_o$,

$$\begin{aligned} \mathbb{E}_{\nu, \gamma} \phi(S, Y) &= \mathbb{E}_{\nu, \gamma} \left(\int_{\mathbb{R}} \phi(S, y) K_\gamma(S, dy) \right) \\ &\leq \mathbb{E}_{\nu, \gamma} \left(\int_{\mathbb{R}} \phi_\alpha(S, y) K_\gamma(S, dy) \right) = \mathbb{E}_{\nu, \gamma} \phi_\alpha(S, Y). \end{aligned}$$

□

Example 9.10 (Fisher's exact test and odds ratios). Suppose that $(X_1, Y_1), \dots, (X_n, Y_n)$ are independent, identically distributed random variables with values in $\{0, 1\} \times \{0, 1\}$ such that all four probabilities

$$p_{xy} := \mathbb{P}(X = x, Y = y), \quad x, y \in \{0, 1\},$$

are strictly positive; here (X, Y) denotes any of the n pairs (X_i, Y_i) . With the parameter $\mathbf{p} := (p_{00}, p_{01}, p_{10}, p_{11})$, the ‘‘correlation’’ between X and Y may be quantified in terms of the odds ratio

$$\rho = \rho(\mathbf{p}) := \frac{p_{11}p_{00}}{p_{01}p_{10}} = \frac{\text{odds}(X = 1 | Y = 1)}{\text{odds}(X = 1 | Y = 0)} = \frac{\text{odds}(Y = 1 | X = 1)}{\text{odds}(Y = 1 | X = 0)}.$$

Suppose we are only interested in ρ . To design good tests or confidence regions for ρ let us rewrite the model as a suitable natural exponential family:

Viewing the X_i and Y_i as random variables on $\Omega := (\{0, 1\} \times \{0, 1\})^n$, equipped with counting measure M , we obtain a statistical experiment with distributions $\mathbb{P}_{\mathbf{p}}$ on Ω such that

$$\begin{aligned} \log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} &= \log(p_{00}^{H_{00}} p_{01}^{H_{01}} p_{10}^{H_{10}} p_{11}^{H_{11}}) \\ &= H_{00} \log p_{00} + H_{01} \log p_{01} + H_{10} \log p_{10} + H_{11} \log p_{11} \end{aligned}$$

with the absolute frequencies

$$H_{xy} := \#\{i \leq n : X_i = x, Y_i = y\}.$$

With the marginal frequencies

$$\begin{aligned} H_{+1} &:= H_{01} + H_{11} = \#\{i \leq n : Y_i = 1\}, \\ H_{1+} &:= H_{10} + H_{11} = \#\{i \leq n : X_i = 1\}, \end{aligned}$$

we may write $H_{00} = n - H_{+1} - H_{1+} + H_{11}$, $H_{01} = H_{+1} - H_{11}$, $H_{10} = H_{1+} - H_{11}$, and this leads to

$$\log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} = H_{+1} \log \frac{p_{01}}{p_{00}} + H_{1+} \log \frac{p_{10}}{p_{00}} + H_{11} \log \frac{p_{11}p_{00}}{p_{01}p_{10}} + n \log p_{00}.$$

With

$$\begin{aligned} \nu = \nu(\mathbf{p}) &:= \left(\log \frac{p_{01}}{p_{00}}, \log \frac{p_{10}}{p_{00}} \right)^\top \in \mathbb{R}^2, \\ \gamma = \gamma(\mathbf{p}) &:= \log \frac{p_{11}p_{00}}{p_{01}p_{10}} = \log \rho \in \mathbb{R}, \end{aligned}$$

we may write

$$\begin{bmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{bmatrix} = (1 + e^{\nu_1} + e^{\nu_2} + e^{\nu_1 + \nu_2 + \gamma})^{-1} \begin{bmatrix} 1 & e^{\nu_1} \\ e^{\nu_2} & e^{\nu_1 + \nu_2 + \gamma} \end{bmatrix}.$$

In particular, for any choice of $\nu \in \mathbb{R}^2$ and $\gamma \in \mathbb{R}$ there exists a probability vector \mathbf{p} such that $\nu = \nu(\mathbf{p})$ and $\gamma = \gamma(\mathbf{p})$. Moreover,

$$\log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} = \nu^\top S + \gamma Y - \kappa(\nu, \gamma)$$

with

$$\begin{aligned} S &:= (H_{+1}, H_{1+})^\top, \\ Y &:= H_{11}, \\ \kappa(\nu, \gamma) &:= -n \log p_{00} = n \log(1 + e^{\nu_1} + e^{\nu_2} + e^{\nu_1 + \nu_2 + \gamma}). \end{aligned}$$

Consequently, if we want to construct tests or confidence regions for ρ , we should concentrate on the conditional distribution of H_{11} , given (H_{+1}, H_{1+}) . For arbitrary integers $s, z \geq 0$, it follows from $H_{+1} = s$ and $H_{1+} = z$ that H_{11} has some value in

$$\{\max(0, s + z - n), \dots, \min(s, z)\}.$$

For any number k in the latter set,

$$\begin{aligned} \mathbb{P}_{\mathbf{p}}(H_{11} = k, H_{+1} = s, H_{1+} = z) &= \mathbb{P}_{\mathbf{p}}(H_{11} = k, H_{01} = s - k, H_{10} = z - k, H_{00} = n - s - z + k) \\ &= \frac{n!}{k!(s-k)!(z-k)!(n-z-s+k)!} p_{11}^k p_{01}^{s-k} p_{10}^{z-k} p_{00}^{n-z-s+k} \\ &= n! p_{01}^s p_{10}^z p_{00}^{n-z-s} \frac{\rho^k}{k!(s-k)!(z-k)!(n-z-s+k)!}. \end{aligned}$$

Consequently,

$$\begin{aligned} \mathbb{P}_{\mathbf{p}}(H_{11} = k \mid H_{+1} = s, H_{1+} = z) &= \frac{\mathbb{P}_{\mathbf{p}}(H_{11} = k, H_{+1} = s, H_{1+} = z)}{\sum_{\ell} \mathbb{P}_{\mathbf{p}}(H_{11} = \ell, H_{+1} = s, H_{1+} = z)} \\ &= C(n, s, z, \rho)^{-1} \frac{\rho^k}{k!(s-k)!(z-k)!(n-z-s+k)!} \end{aligned}$$

with

$$C(n, s, z, \rho) := \sum_{\ell=\max(0, s+z-n)}^{\min(s, z)} \frac{\rho^\ell}{\ell!(s-\ell)!(z-\ell)!(n-z-s+\ell)!}.$$

Alternatively one may write

$$\mathbb{P}_{\mathbf{p}}(H_{11} = k \mid H_{+1} = s, H_{1+} = z) = \tilde{C}(n, s, z, \rho)^{-1} \text{Hyp}_{n, s, z}(\{k\}) \rho^k$$

with the hypergeometric distribution $\text{Hyp}_{n, s, z}$ and

$$\tilde{C}(n, s, z, \rho) := \sum_{\ell=\max(0, s+z-n)}^{\min(s, z)} \text{Hyp}_{n, s, z}(\{\ell\}) \rho^\ell.$$

In particular, if $\rho = 1$, the conditional distribution of H_{11} , given $H_{+1} = s$ and $H_{1+} = z$, equals $\text{Hyp}_{n, s, z}$.

Example 9.11 (McNemar's test). Traditionally, McNemar's test is described in the context of two-by-two tables as in the previous example. But it may be transferred to a more general setting: Let X_1, X_2, \dots, X_n be independent and identically distributed random variables with values in

a finite set $\mathcal{X} := \{x_1, x_2, \dots, x_K\}$ with $K \geq 3$ different elements. The parameter vector $\mathbf{p} = (p_k)_{k=1}^K$ with $p_k := \mathbb{P}(X = x_k) > 0$ is unknown. Suppose we are mainly interested in the ratio

$$\rho = \rho(\mathbf{p}) := \frac{p_1}{p_2}.$$

To this end we consider the X_i as random variables on the finite set $\Omega := \mathcal{X}^n$, equipped with counting measure M . This leads to a statistical experiment with distributions $\mathbb{P}_{\mathbf{p}}$ on Ω given by

$$\log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} = \log \prod_{j=1}^K p_j^{H_j} = \sum_{j=1}^K H_j \log(p_j)$$

with the absolute frequencies

$$H_j := \#\{i \leq n : X_i = x_j\}.$$

Since $H_K = n - \sum_{j < K} H_j$, we may rewrite this as

$$\begin{aligned} \log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} &= \sum_{j=1}^{K-1} H_j \log(p_j/p_K) + n \log(p_K) \\ &= H_1 \log \frac{p_1}{p_2} + (H_1 + H_2) \log \frac{p_2}{p_K} + \sum_{2 < j < K} H_j \log \frac{p_j}{p_K} + n \log(p_K). \end{aligned}$$

With

$$\begin{aligned} \nu &= \nu(\mathbf{p}) := \left(\log \frac{p_{\ell+1}}{p_K} \right)_{\ell=1}^{K-2} \in \mathbb{R}^{K-2}, \\ \gamma &= \gamma(\mathbf{p}) := \log \frac{p_1}{p_2} = \log \rho \end{aligned}$$

we may write

$$\mathbf{p} = \frac{(e^{\gamma+\nu_1}, e^{\nu_1}, \dots, e^{\nu_{K-2}}, 1)^\top}{e^{\gamma+\nu_1} + e^{\nu_1} + \dots + e^{\nu_{K-2}} + 1}.$$

In particular, for any choice of $\nu \in \mathbb{R}^{K-2}$ and $\gamma \in \mathbb{R}$ there exists a probability vector \mathbf{p} such that $\nu = \nu(\mathbf{p})$ and $\gamma = \gamma(\mathbf{p})$. Moreover,

$$\log \frac{d\mathbb{P}_{\mathbf{p}}}{dM} = \nu^\top S + \gamma Y - \kappa(\nu, \gamma)$$

with

$$\begin{aligned} S_1 &:= H_1 + H_2, \\ S_\ell &:= H_{\ell+1} \quad \text{for } 2 \leq \ell \leq K-2, \\ Y &:= H_1, \\ \kappa(\nu, \gamma) &:= -n \log(p_K) = n \log(e^{\gamma+\nu_1} + e^{\nu_1} + \dots + e^{\nu_{K-2}} + 1). \end{aligned}$$

Consequently, if we want to construct tests or confidence regions for ρ , we should concentrate on the conditional distribution of H_1 , given $(H_1 + H_2, H_3, \dots, H_K)$. For arbitrary integers

$m, s_3, \dots, s_K \geq 0$ with $m + \sum_{j=3}^K s_j = n$ and $k \in \{0, \dots, m\}$,

$$\begin{aligned} & \mathbb{P}_{\mathbf{p}}(H_1 = k, H_1 + H_2 = m, H_j = s_j \text{ for } j \geq 3) \\ &= \mathbb{P}_{\mathbf{p}}(H_1 = k, H_2 = m - k, H_j = s_j \text{ for } j \geq 3) \\ &= \frac{n!}{k!(m-k)! \prod_{j \geq 3} s_j!} p_1^k p_2^{m-k} \prod_{j \geq 3} p_j^{s_j} \\ &= \binom{m}{k} \pi^k (1 - \pi)^{m-k} \frac{n!}{m! \prod_{j \geq 3} s_j!} (p_1 + p_2)^m \prod_{j \geq 3} p_j^{s_j} \end{aligned}$$

with

$$\pi = \pi(\mathbf{p}) := \frac{p_1}{p_1 + p_2} = \frac{\rho}{1 + \rho} \in (0, 1).$$

Consequently, the conditional distribution of H_1 , given $H_1 + H_2, H_3, \dots, H_K$, equals

$$\text{Bin}(H_1 + H_2, \pi) = \text{Bin}\left(H_1 + H_2, \frac{\rho}{1 + \rho}\right).$$

Any test or confidence region for π may be translated into a test or a confidence region for ρ via the inverse transformation

$$\rho = \frac{\pi}{1 - \pi}.$$

Example 9.12 (Comparing two Poisson parameters). Suppose we observe independent random variables $X \sim \text{Poiss}(\lambda)$ and $Y \sim \text{Poiss}(\mu)$ with unknown parameters $\lambda, \mu > 0$. Suppose further that we are mainly interested in the ratio $\rho := \lambda/\mu$. With $\Omega := \mathbb{N}_0 \times \mathbb{N}_0$, $X(\omega) := \omega_1$, $Y(\omega) := \omega_2$ and counting measure M on Ω , this situation corresponds to a statistical model with distributions $\mathbb{P}_{\lambda, \mu}$ given by

$$\frac{d\mathbb{P}_{\lambda, \mu}}{dM} = e^{-\lambda} \frac{\lambda^X}{X!} e^{-\mu} \frac{\mu^Y}{Y!} = \frac{\exp(X \log \lambda + Y \log \mu - \lambda - \mu)}{X! Y!}.$$

Replacing M with the measure M_o given by $M_o(\{\omega\}) := (\omega_1! \omega_2!)^{-1}$ we may write

$$\begin{aligned} \log \frac{d\mathbb{P}_{\lambda, \mu}}{dM_o} &= X \log \lambda + Y \log \mu - (\lambda + \mu) \\ &= X \log(\lambda/\mu) + (X + Y) \log \mu - (\lambda + \mu) \\ &= \nu(X + Y) + \gamma X - \kappa(\nu, \gamma) \end{aligned}$$

with

$$\begin{aligned} \nu &:= \log \mu, \\ \gamma &:= \log(\lambda/\mu) = \log \rho, \\ \kappa(\nu, \gamma) &:= \lambda + \mu = \exp(2\nu + \gamma). \end{aligned}$$

Consequently, for inference about ρ we should analyze the conditional distribution of X , given $X + Y$. But for arbitrary integers $m \geq 0$ and $k \in \{0, \dots, m\}$,

$$\begin{aligned} \mathbb{P}_{\lambda, \mu}(X = k, X + Y = m) &= \mathbb{P}_{\lambda, \mu}(X = k) \mathbb{P}_{\lambda, \mu}(Y = m - k) \\ &= e^{-(\lambda + \mu)} \frac{\lambda^k \mu^{m-k}}{k!(m-k)!} \\ &= \binom{m}{k} \pi^k (1 - \pi)^{m-k} e^{-(\lambda + \mu)} \frac{(\lambda + \mu)^m}{m!} \\ &= \text{Bin}_{m, \pi}(\{k\}) \text{Pois}_{\lambda + \mu}(\{m\}) \end{aligned}$$

with

$$\pi := \frac{\lambda}{\lambda + \mu} = \frac{\rho}{1 + \rho} \in (0, 1).$$

This shows that the conditional distribution of X , given $X + Y$, is equal to

$$\text{Bin}(X + Y, \pi) = \text{Bin}\left(X + Y, \frac{\rho}{1 + \rho}\right).$$

Hence we may construct tests and confidence regions for π , and these translate into tests and confidence regions for ρ via the inverse transformation $\rho = \pi/(1 - \pi)$.

Avoiding conditional distributions

In the previous examples we computed the conditional distribution of Y , given S , explicitly. In various settings this step can be avoided by means of the following result:

Lemma 9.13 (Basú). *Let $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ be a statistical experiment, and let $S : (\Omega, \mathcal{A}) \rightarrow (\Omega', \mathcal{A}')$ be a sufficient statistic for \mathcal{E} such that the experiment $\mathcal{E}^S = (\Omega', \mathcal{A}', (\mathbb{P}_\theta^S)_{\theta \in \Theta})$ is boundedly complete. If $V : (\Omega, \mathcal{A}) \rightarrow (\Omega'', \mathcal{A}'')$ is a measurable mapping such that its distributions \mathbb{P}_θ^V , $\theta \in \Theta$, are identical, then S and V are stochastically independent under each measure \mathbb{P}_θ , $\theta \in \Theta$.*

Proof of Lemma 9.13. By sufficiency of S there exists a stochastic kernel K from (Ω, \mathcal{A}) to (Ω', \mathcal{A}') describing the conditional distribution of $X \sim \mathbb{P}_\theta$, given S , simultaneously for all $\theta \in \Theta$. That means, for arbitrary $A' \in \mathcal{A}'$, $A'' \in \mathcal{A}''$ and $\theta \in \Theta$,

$$\begin{aligned} \mathbb{P}_\theta(S \in A', V \in A'') &= \mathbb{E}_\theta(\mathbb{E}_\theta(1_{A'}(S)1_{A''}(V) | S)) \\ &= \mathbb{E}_\theta(1_{A'}(S) \mathbb{P}_\theta(V \in A'' | S)) \\ &= \mathbb{E}_\theta(1_{A'}(S)K(S, \{V \in A''\})). \end{aligned}$$

Setting $A' = \Omega'$, we obtain the equation

$$\mathbb{P}_\theta(V \in A'') = \mathbb{E}_\theta K(S, \{V \in A''\}).$$

By our assumption on V , the left hand side does not depend on $\theta \in \Theta$, and we denote this number with $P(V \in A'')$. Hence $f(s) := K(s, \{V \in A''\}) - P(V \in A'')$ defines a bounded measurable function on (Ω', \mathcal{A}'') such that $\mathbb{E}_\theta f(S) = 0$ for all $\theta \in \Theta$. By completeness of \mathcal{E}^S ,

$$\mathbb{P}_\theta(f(S) \neq 0) = \mathbb{P}_\theta(K(S, \{V \in A''\}) \neq P(V \in A'')) = 0 \quad \text{for all } \theta \in \Theta.$$

Hence for arbitrary $A' \in \mathcal{A}'$, $A'' \in \mathcal{A}''$ and $\theta \in \Theta$,

$$\begin{aligned} \mathbb{P}_\theta(S \in A', V \in A'') &= \mathbb{E}_\theta(1_{A'}(S)K(S, \{V \in A''\})) \\ &= \mathbb{E}_\theta(1_{A'}(S)P(V \in A'')) \\ &= \mathbb{P}_\theta(S \in A')P(V \in A''), \end{aligned}$$

which proves stochastic independence of S and V under \mathbb{P}_θ . □

Application to exponential families. Let $\mathcal{E} = (\Omega, \mathcal{A}, (\mathbb{P}_\theta)_{\theta \in \Theta})$ be a natural exponential family as described at the beginning of this section with open and convex parameter space $\Theta = N \times \Gamma \subset \mathbb{R}^d \times \mathbb{R}$ and sufficient statistic $(S, Y) \in \mathbb{R}^d \times \mathbb{R}$. Writing $\theta \in \Theta$ as $\theta = (\nu, \gamma)$, we know from Neyman's theory how to construct optimal unbiased level- α tests of the null hypotheses “ $\gamma \leq \gamma_o$ ” or “ $\gamma \geq \gamma_o$ ” or “ $\gamma = \gamma_o$ ” for any given value $\gamma_o \in \Gamma$.

Note that for the restricted experiment $\mathcal{E}_o = (\Omega, \mathcal{A}, (\mathbb{P}_{\nu, \gamma_o})_{\nu \in N})$ the statistic S is sufficient, and the family $\mathcal{E}_o^S = (\mathbb{R}^d, \text{Borel}(\mathbb{R}^d), (\mathbb{P}_{\nu, \gamma_o}^S)_{\nu \in N})$ is complete. Suppose we can identify a real-valued statistic

$$V = f(S, Y)$$

with the following two properties:

- $V = f(S, Y)$ is strictly increasing in Y almost surely,
- the distribution of V under $\mathbb{P}_{\nu, \gamma_o}$ does not depend on $\nu \in N$.

Then optimal unbiased level- α tests of the null hypotheses above may be constructed in terms of V and its unconditional distribution under $\mathbb{P}_{\nu, \gamma_o}$, where $\nu \in N$ is arbitrary.

For instance if V has continuous distribution function F_{γ_o} in case of $\gamma = \gamma_o$, then the right-sided p-value

$$1 - F_{\gamma_o}(V)$$

yields an optimal unbiased level- α test of “ $\gamma \leq \gamma_o$ ”, whereas the left-sided p-value

$$F_{\gamma_o}(V)$$

is optimal for the null hypothesis “ $\gamma \geq \gamma_o$ ”.

Example 9.14 (Student's t -test for a Gaussian mean). As in Example 9.3 we consider the statistical experiment $\mathcal{E} = (\mathbb{R}^n, \text{Borel}(\mathbb{R}^n), (\mathcal{N}(\mu, \sigma^2)^{\otimes n})_{(\mu, \sigma) \in \mathbb{R} \times (0, \infty)})$. Suppose we want to construct optimal unbiased level- α tests of “ $\mu \leq \mu_o$ ” or “ $\mu \geq \mu_o$ ” or “ $\mu = \mu_o$ ”, where μ_o is a given fixed number. In all three cases we have to deal with the nuisance parameter σ . With M denoting Lebesgue measure on \mathbb{R}^n times $(2\pi)^{-n/2}$ and $X_i(\omega) := \omega_i$, we may write

$$\begin{aligned} \log \frac{d\mathbb{P}_{\mu, \sigma}}{dM} &= - \sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2} - n \log \sigma \\ &= - \sum_{i=1}^n \frac{(X_i - \mu_o)^2}{2\sigma^2} + \frac{n(\mu - \mu_o)(\bar{X} - \mu_o)}{\sigma^2} - n \log \sigma \\ &= \nu S + \gamma Y - \kappa(\nu, \gamma), \end{aligned}$$

where

$$\begin{aligned}\nu &:= \frac{-1}{2\sigma^2} \in (-\infty, 0), \\ \gamma = \gamma(\mu_o) &:= \frac{\sqrt{n}(\mu - \mu_o)}{\sigma^2} \in \mathbb{R}, \\ \kappa(\nu, \gamma) &:= n \log \sigma, \\ S = S(\mu_o) &:= \sum_{i=1}^n (X_i - \mu_o)^2, \\ Y = Y(\mu_o) &:= \sqrt{n}(\bar{X} - \mu_o).\end{aligned}$$

Hence an unbiased test of “ $\mu \leq \mu_o$ ” (or of “ $\mu = \mu_o$ ”) may be identified with an unbiased test of “ $\gamma \leq 0$ ” (or of “ $\gamma = 0$ ”).

Instead of determining the conditional distribution of Y , given S , under $\mathbb{P}_{\mu_o, \sigma}$ for some $\sigma > 0$ directly, we apply Basú’s lemma and recall student’s method from introductory statistics courses: With the sample standard deviation

$$S_X := \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

it is well-known that

$$V = V(\mu_o) := \frac{\sqrt{n}(\bar{X} - \mu_o)}{S_X} \sim t_{n-1}$$

whenever $\mu = \mu_o$, irrespective of $\sigma > 0$. But

$$\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i - \mu_o)^2 - n(\bar{X} - \mu_o)^2 = S - Y^2,$$

so

$$V = \frac{\sqrt{n-1}Y}{\sqrt{S - Y^2}},$$

which is monotone increasing in Y almost surely. (Note that $S > 0$ and $Y^2 < S$ almost surely.) Hence with F_{n-1} denoting the distribution function of t_{n-1} , the left-sided p-value

$$F_{n-1}(V(\mu_o))$$

yields optimal unbiased tests of “ $\mu \geq \mu_o$ ”, the right-sided p-value

$$1 - F_{n-1}(V(\mu_o))$$

yields optimal unbiased tests of “ $\mu \leq \mu_o$ ”, and the traditional $(1 - \alpha)$ -confidence interval

$$\left[\bar{X} \pm \frac{S_X}{\sqrt{n}} t_{n-1; 1-\alpha/2} \right]$$

for μ is based on optimal unbiased level- α tests of one-point hypotheses “ $\mu = \mu_o$ ”, $\mu_o \in \mathbb{R}$.

Example 9.15 (Comparing two Gamma scale parameters). Suppose we observe independent random variables

$$Y_1 \sim \text{Gamma}(a_1, \beta_1) \quad \text{and} \quad Y_2 \sim \text{Gamma}(a_2, \beta_2)$$

with given shape parameters $a_1, a_2 > 0$ and unknown scale parameters $\beta_1, \beta_2 > 0$. Suppose we are mainly interested in the ratio

$$\rho = \rho(\boldsymbol{\beta}) := \beta_1/\beta_2,$$

where $\boldsymbol{\beta} = (\beta_1, \beta_2) \in (0, \infty) \times (0, \infty)$. With $\Omega := (0, \infty)^2$, Lebesgue measure M on Ω and $Y_i(\omega) := \omega_i$, this corresponds to the statistical model with distributions $\mathbb{P}_{\boldsymbol{\beta}}$ given by

$$\frac{d\mathbb{P}_{\boldsymbol{\beta}}}{dM} = \frac{(Y_1/\beta_1)^{a_1-1} \exp(-Y_1/\beta_1)}{\Gamma(a_1)\beta_1} \frac{(Y_2/\beta_2)^{a_2-1} \exp(-Y_2/\beta_2)}{\Gamma(a_2)\beta_2}.$$

With the modified measure M_o given by

$$\frac{dM_o}{dM} := \frac{Y_1^{a_1-1} Y_2^{a_2-1}}{\Gamma(a_1)\Gamma(a_2)}$$

we may write

$$\log \frac{d\mathbb{P}_{\boldsymbol{\beta}}}{dM_o} = \frac{-1}{\beta_1} Y_1 + \frac{-1}{\beta_2} Y_2 - a_1 \log \beta_1 - a_2 \log \beta_2.$$

For a hypothetical value ρ_o of ρ we may rewrite the log-likelihood as

$$\begin{aligned} \log \frac{d\mathbb{P}_{\boldsymbol{\beta}}}{dM_o} &= \frac{\rho - \rho_o}{\beta_2 \rho} Y_1/\rho_o + \frac{-1}{\beta_2} (Y_1/\rho_o + Y_2) - a_1 \log \beta_1 - a_2 \log \beta_2 \\ &= \nu S + \gamma Y - \kappa(\nu, \gamma) \end{aligned}$$

with

$$\begin{aligned} \nu &= \nu(\boldsymbol{\beta}) := \frac{-1}{\beta_2} \in (-\infty, 0), \\ \gamma &= \gamma(\boldsymbol{\beta}, \rho_o) := \frac{\rho - \rho_o}{\beta_2 \rho} \in \mathbb{R}, \\ S &= S(\rho_o) := Y_1/\rho_o + Y_2, \\ Y &= Y(\rho_o) := Y_1/\rho_o \end{aligned}$$

and the normalization constant $\kappa(\nu, \gamma) = a_1 \log \beta_1 + a_2 \log \beta_2$. Note that (ν, γ) lies in the open convex set

$$\{(\nu, \gamma) \in (-\infty, 0) \times \mathbb{R} : \gamma < -\nu\},$$

which contains $(-\infty, 0) \times \{0\}$. Hence optimal unbiased tests of the null hypothesis “ $\rho \leq \rho_o$ ” or “ $\rho \geq \rho_o$ ” or “ $\rho = \rho_o$ ” may be viewed as optimal unbiased tests of the null hypothesis “ $\gamma \leq 0$ ” or “ $\gamma \geq 0$ ” or “ $\gamma = 0$ ” and could be constructed with the conditional distribution of Y , given S , under any distribution $\mathbb{P}_{\beta_2 \rho_o, \beta_2}$.

In case of $\rho = \rho_o$, the distribution of $(Y_1/\rho_o, Y_2)$ coincides with the distribution of $\beta_2(Z_1, Z_2)$ with independent random variables

$$Z_1 \sim \text{Gamma}(a_1, 1) \quad \text{and} \quad Z_2 \sim \text{Gamma}(a_2, 1).$$

Then the test statistic

$$V(\rho_o) := \frac{Y(\rho_o)}{S(\rho_o)} = \frac{Y_1/\rho_o}{Y_1/\rho_o + Y_2}$$

has the same distribution as

$$\frac{Z_1}{Z_1 + Z_2} \sim \text{Beta}(a_1, a_2),$$

irrespective of $\beta_2 = -1/\nu$. Thus by Basu's lemma or a direct argument, $V(\rho_o)$ and $S(\rho_o)$ are stochastically independent in case of $\rho = \rho_o$, and optimal tests are obtained by comparing $V(\rho_o)$ with $\text{Beta}(a_1, a_2)$.

Specifically, if we denote the u -quantile of $\text{Beta}(a_1, a_2)$ with $\text{Beta}_{a_1, a_2; u}$, then an optimal unbiased level- α test of " $\rho \geq \rho_o$ " rejects this null hypothesis if

$$V(\rho_o) \leq \text{Beta}_{a_1, a_2; \alpha}.$$

This leads to the $(1 - \alpha)$ -confidence region

$$\begin{aligned} C_\alpha(Y_1, Y_2) &:= \{ \rho_o > 0 : V(\rho_o) > \text{Beta}_{a_1, a_2; \alpha} \} \\ &= \left\{ \rho_o > 0 : \frac{Y_1}{Y_1 + \rho_o Y_2} > \text{Beta}_{a_1, a_2; \alpha} \right\} \\ &= \left(0, \frac{Y_1}{Y_2} (B_{a_1, a_2; \alpha}^{-1} - 1) \right). \end{aligned}$$

for ρ . Analogously, an optimal unbiased level- α test of " $\rho \leq \rho_o$ " rejects this null hypothesis if

$$V(\rho_o) \geq \text{Beta}_{a_1, a_2; 1-\alpha},$$

and this leads to the $(1 - \alpha)$ -confidence region

$$C_\alpha(Y_1, Y_2) := \left(\frac{Y_1}{Y_2} (B_{a_1, a_2; 1-\alpha}^{-1} - 1), \infty \right).$$

for ρ .

Chapter 10

Some Asymptotics

10.1 Testing, Total Variation and Hellinger Distances

Let P, Q be two probability distributions on a measurable space (Ω, \mathcal{A}) . In this section we introduce statistically meaningful measures of distance between P and Q and establish connections between them.

Testing affinity and distance. Let X be a random variable with unknown distribution in $\{P, Q\}$. Now we consider a statistical test $\varphi : \Omega \rightarrow [0, 1]$ and interpret $\varphi(X)$ as the probability of claiming that $X \sim Q$, whereas $1 - \varphi(X)$ is the probability of claiming that $X \sim P$. The quality of φ could be measured by the risk

$$\mathbb{P}(\text{error of 1st kind}) + \mathbb{P}(\text{error of 2nd kind}) = \int \varphi dP + \int (1 - \varphi) dQ.$$

Hence a natural measure of similarity between P and Q is given by the *testing affinity*

$$\eta_T(P, Q) := \inf_{\text{tests } \varphi} \left(\int \varphi dP + \int (1 - \varphi) dQ \right).$$

As shown in the next lemma, the latter infimum is always a minimum.

Lemma 10.1 (Testing affinity and distance). *Suppose that P and Q have densities f and g , respectively, with respect to some measure M on (Ω, \mathcal{A}) . Then*

$$\eta_T(P, Q) = \int \min(f, g) dM.$$

The minimum is attained by any test φ such that $\varphi = 0$ on $\{f > g\}$ and $\varphi = 1$ on $\{f < g\}$.

There exists always a measure M dominating both P and Q , and the testing affinity $\eta_T(P, Q)$ does not depend on the choice of M . The related quantity

$$D_T(P, Q) := 1 - \int \min(f, g) dM = \frac{1}{2} \int |f - g| dM$$

is the so-called testing distance between P and Q . It defines a metric on the space of probability measures on (Ω, M) with values in $[0, 1]$.

Proof of Lemma 10.1. Note first that for any test φ ,

$$\begin{aligned} \int \varphi dP + \int (1 - \varphi) dQ &= \int (\varphi(\omega)f(\omega) + (1 - \varphi(\omega))g(\omega)) M(d\omega) \\ &\geq \int \min(f(\omega), g(\omega)) M(d\omega) \end{aligned}$$

with equality if

$$\varphi(\omega) = \begin{cases} 0 & \text{if } f(\omega) > g(\omega), \\ 1 & \text{if } f(\omega) < g(\omega). \end{cases}$$

This proves already the first part of the lemma.

As to the second part, consider the particular measure $M_o := P + Q$. Then $M_o(\Omega) = 2$ and $P(A) = Q(A) = 0$ whenever $M_o(A) = 0$. Hence by the theorem of Radon–Nikodym there exist densities f_o and g_o of P and Q , respectively, with respect to M_o .

Now let M be an arbitrary measure such that $f := dP/dM$ and $g := dQ/dM$ exist. Then M_o has density $h := f + g$ with respect to M . Defining $f_o := f/h$ and $g_o := g/h$ with the convention $0/0 := 0$, we have

$$\int \min(f_o, g_o) dM_o = \int \min(f/h, g/h)h dM = \int \min(f, g) dM.$$

This shows that the infimum in $\eta_T(P, Q)$ is always a minimum, and its value does not depend on the choice of M .

Since $f, g \geq \min(f, g) \geq 0$, it is clear that $0 \leq \int \min(f, g) dM \leq 1$. Moreover,

$$\begin{aligned} 1 - \eta_T(P, Q) &= \int \left(\frac{f + g}{2} - \min(f, g) \right) dM \\ &= \frac{1}{2} \int (f + g - 2 \min(f, g)) dM \\ &= \frac{1}{2} \int (\max(f, g) - \min(f, g)) dM \\ &= \frac{1}{2} \int |f - g| dM. \end{aligned}$$

Obviously this equals 0 if, and only if, $f = g$ M -almost everywhere, which is equivalent to $P \equiv Q$.

Finally, if P, Q, R are probability distributions on (Ω, \mathcal{A}) , we may assume that there exists a finite measure M on (Ω, \mathcal{A}) such that $f = dP/dM$, $g = dQ/dM$ and $h = dR/dM$ exist. For instance, $M := P + Q + R$ would do the job. But then by the triangle inequality for the norm $\|\cdot\|_{M,1}$ in $L^1(M)$,

$$D_T(P, R) = \frac{\|f - h\|_{M,1}}{2} \leq \frac{\|f - g\|_{M,1} + \|g - h\|_{M,1}}{2} = D_T(P, Q) + D_T(Q, R).$$

□

Total variation distance. Another measure of distance is given by the total variation distance

$$D_{TV}(P, Q) := \sup_{A \in \mathcal{A}} |P(A) - Q(A)|.$$

One could easily verify directly that this defines a metric on the space of probability measures on (Ω, \mathcal{A}) . But the next lemma shows that $D_{TV} = D_T$.

Lemma 10.2. For arbitrary probability distributions P and Q on (Ω, \mathcal{A}) ,

$$D_{TV}(P, Q) = \max_{A \in \mathcal{A}} (P(A) - Q(A)) = \max_{B \in \mathcal{A}} (Q(B) - P(B)) = D_T(P, Q).$$

If P and Q have densities f and g as in Lemma 10.1, then the latter two maxima are attained for arbitrary sets A, B such that $\{f > g\} \subset A \subset \{f \geq g\}$ and $\{f < g\} \subset B \subset \{f \leq g\}$.

Proof of Lemma 10.2. Since $P(A) - Q(A) = Q(\Omega \setminus A) - P(\Omega \setminus A)$,

$$D_{TV}(P, Q) = \sup_{A \in \mathcal{A}} (P(A) - Q(A)) = \sup_{B \in \mathcal{A}} (Q(B) - P(B)).$$

In case of P and Q having densities f and g , respectively, with respect to some measure M on (Ω, \mathcal{A}) ,

$$P(A) - Q(A) = \int 1_A(f - g) dM \leq \int (f - g)^+ dM$$

with equality if $\{f > g\} \subset A \subset \{f \geq g\}$. Analogously,

$$Q(B) - P(B) \leq \int (g - f)^+ dM = \int (f - g)^- dM$$

with equality if $\{f < g\} \subset B \subset \{f \leq g\}$. Consequently,

$$D_{TV}(P, Q) = \frac{1}{2} \int ((f - g)^+ + (f - g)^-) dM = \frac{1}{2} \int |f - g| dM = D_T(P, Q). \quad \square$$

Hellinger affinity and distance. In many situations it turns out that $D_T(P, Q)$ is difficult to compute explicitly. As we shall see later, interesting proxys are given by the

$$\text{Hellinger affinity } \eta_H(P, Q) := \int \sqrt{fg} dM,$$

and

$$\text{Hellinger distance } D_H(P, Q) := \sqrt{\frac{1}{2} \int (\sqrt{f} - \sqrt{g})^2 dM},$$

where M is some measure on (Ω, \mathcal{A}) such that $f := dP/dM$ and $g := dQ/dM$ exist. As in case of $\eta_T(P, Q)$ and $D_T(P, Q)$, the choice of M is irrelevant. Note also that

$$D_H(P, Q)^2 = \frac{1}{2} \int (f + g - 2\sqrt{fg}) dM = 1 - \eta_H(P, Q).$$

The following lemma shows that testing and Hellinger distance induce the same topology on the space of probability distributions on (Ω, \mathcal{A}) .

Lemma 10.3 (Relationships between testing and Hellinger distance).

$$1 - \sqrt{1 - \eta_H^2} \leq \eta_T \leq \eta_H$$

and

$$D_H^2 \leq D_T \leq D_H \sqrt{2 - D_H^2}.$$

Proof of Lemma 10.3. With explicit densities $f := dP/dM$ and $g := dQ/dM$ it follows from $\min(a, b) \leq \sqrt{ab}$ for real numbers $a, b \geq 0$ that

$$\eta_T(P, Q) = \int \min(f, g) dM \leq \int \sqrt{fg} dM = \eta_H(P, Q).$$

In particular,

$$D_T(P, Q) = 1 - \eta_T(P, Q) \geq 1 - \eta_H(P, Q) = D_H(P, Q)^2.$$

As to the other bounds, it follows from $a - b = (\sqrt{a} - \sqrt{b})(\sqrt{a} + \sqrt{b})$ for real numbers $a, b \geq 0$ and the Cauchy–Schwarz inequality that

$$\begin{aligned} 1 - \eta_T(P, Q) = D_T(P, Q) &= \frac{1}{2} \int |\sqrt{f} - \sqrt{g}| |\sqrt{f} + \sqrt{g}| dM \\ &\leq \frac{1}{2} \sqrt{\int (f + g - 2\sqrt{fg}) dM \int (f + g + 2\sqrt{fg}) dM} \\ &= \sqrt{(1 - \eta_H(P, Q))(1 + \eta_H(P, Q))} \\ &= \begin{cases} \sqrt{1 - \eta_H(P, Q)^2}, \\ D_H(P, Q) \sqrt{2 - D_H(P, Q)^2}. \end{cases} \end{aligned}$$

This proves that

$$D_T \leq D_H \sqrt{2 - D_H^2} \quad \text{and} \quad 1 - \eta_T \leq \sqrt{1 - \eta_H^2},$$

where the latter inequality is equivalent to $\eta_T \geq 1 - \sqrt{1 - \eta_H^2}$. \square

Remark 10.4. As mentioned already, the formulae for η_T , D_T , η_H and D_H are independent of the choice of the dominating measure M . Thus some authors write symbolically

$$\begin{aligned} \eta_T(P, Q) &= \int \min(dP, dQ), \\ D_T(P, Q) &= \frac{1}{2} \int |dP - dQ|, \\ \eta_H(P, Q) &= \int \sqrt{dP dQ}, \\ D_H(P, Q) &= \sqrt{\frac{1}{2} \int (\sqrt{dP} - \sqrt{dQ})^2}. \end{aligned}$$

Remark 10.5. Let (Ω', \mathcal{A}') be a second measurable space, and let $\tau : \Omega \rightarrow \Omega'$ be a bijection such that both τ and τ^{-1} are measurable. Then

$$\kappa(P, Q) = \kappa(P^\tau, Q^\tau) \quad \text{for } \kappa = \eta_T, D_T, \eta_H, D_H.$$

The proof of this claim is left to the reader as an exercise; see also Exercise 2.19.

Example 10.6 (Testing and Hellinger distance for univariate Gaussian shift). For real numbers μ_1, μ_2 and $\sigma > 0$,

$$\begin{aligned}\eta_T(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) &= 2\Phi\left(-\frac{|\mu_1 - \mu_2|}{2\sigma}\right), \\ \eta_H(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) &= \exp\left(-\frac{(\mu_1 - \mu_2)^2}{8\sigma^2}\right),\end{aligned}$$

where Φ is the standard Gaussian distribution function. Hence

$$\begin{aligned}D_T(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) &= 2\Phi\left(\frac{|\mu_1 - \mu_2|}{2\sigma}\right) - 1, \\ D_H(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) &= \sqrt{1 - \exp\left(-\frac{(\mu_1 - \mu_2)^2}{8\sigma^2}\right)}.\end{aligned}$$

To verify these formulae, we first apply Remark 10.5 to the bijection $\tau : \mathbb{R} \rightarrow \mathbb{R}$ with $\tau(x) := \sigma^{-1}(x - \min(\mu_1, \mu_2))$. Then it suffices to verify the asserted formulae for the testing and Hellinger affinity in case of $P = \mathcal{N}(0, 1)$ and $Q = \mathcal{N}(\mu, 1)$, where

$$\mu := \frac{|\mu_1 - \mu_2|}{\sigma}.$$

But with $C = (2\pi)^{-1/2}$ and $\phi(x) := C \exp(-x^2/2)$,

$$\begin{aligned}\eta_T(\mathcal{N}(0, 1), \mathcal{N}(\mu, 1)) &= \int \min(\phi(x), \phi(x - \mu)) dx \\ &= C \int \exp\left(-\frac{\max(x^2, (x - \mu)^2)}{2}\right) dx \\ &= 2C \int_{-\infty}^{\mu/2} \exp\left(-\frac{(x - \mu)^2}{2}\right) dx \\ &= 2 \int_{-\infty}^{\mu/2} \phi(x - \mu) dx \\ &= 2\Phi(-\mu/2).\end{aligned}$$

Moreover,

$$\begin{aligned}\eta_H(\mathcal{N}(0, 1), \mathcal{N}(\mu, 1)) &= \int \sqrt{\phi(x)\phi(x - \mu)} dx \\ &= C \int \exp\left(-\frac{x^2 + (x - \mu)^2}{4}\right) dx \\ &= C \int \exp\left(-\frac{x^2 - x\mu + \mu^2/2}{2}\right) dx \\ &= C \int \exp\left(-\frac{(x - \mu/2)^2 + \mu^2/4}{2}\right) dx \\ &= \exp(-\mu^2/8) \int \phi(x - \mu/2) dx \\ &= \exp(-\mu^2/8).\end{aligned}$$

Exercise 10.7 (Hellinger distance for multivariate Gaussian shift). For any dimension $d \geq 1$, consider arbitrary vectors $\mu_1, \mu_2 \in \mathbb{R}^d$ and a symmetric, positive definite matrix $\Sigma \in \mathbb{R}^{d \times d}$. Show

that

$$\begin{aligned}\eta_T(\mathcal{N}_d(\mu_1, \Sigma), \mathcal{N}_d(\mu_2, \Sigma)) &= 2\Phi\left(-\sqrt{(\mu_2 - \mu_1)^\top \Sigma^{-1}(\mu_2 - \mu_1)}/2\right), \\ \eta_H(\mathcal{N}_d(\mu_1, \Sigma), \mathcal{N}_d(\mu_2, \Sigma)) &= \exp(-(\mu_2 - \mu_1)^\top \Sigma^{-1}(\mu_2 - \mu_1)/8).\end{aligned}$$

Hint: Consider the transformation $\tau(x) := \Sigma^{-1/2}(x - \mu_1)$ of $x \in \mathbb{R}^d$.

Exercise 10.8 (More about the relation between testing and Hellinger distance). **(a)** Show that the inequalities

$$D_H^2 \leq D_T \leq D_H \sqrt{2 - D_H^2}$$

are equivalent to

$$\sqrt{1 - \sqrt{1 - D_T^2}} \leq D_H \leq \sqrt{D_T}.$$

(b) Visualize these two pairs of bounds graphically.

(c) Construct for any $\gamma \in [0, 1]$ two distributions P_γ and Q_γ on a suitable sample space (Ω, \mathcal{A}) such that $D_T(P_\gamma, Q_\gamma) = D_H^2(P_\gamma, Q_\gamma) = \gamma$.

Remark 10.9 (Product measures). For $j = 1, 2$, let P_j and Q_j be probability measures on a measurable space $(\Omega_j, \mathcal{A}_j)$. Then

$$\eta_H(P_1 \otimes P_2, Q_1 \otimes Q_2) = \eta_H(P_1, Q_1)\eta_H(P_2, Q_2).$$

For if $dP_j/dM_j = f_j$ and $dQ_j/dM_j = g_j$, then by Fubini's theorem,

$$\begin{aligned}\eta_H(P_1 \otimes P_2, Q_1 \otimes Q_2) &= \int_{\Omega_1 \times \Omega_2} \sqrt{f_1(\omega_1)f_2(\omega_2)g_1(\omega_1)g_2(\omega_2)} M_1 \otimes M_2(d(\omega_1, \omega_2)) \\ &= \int_{\Omega_1} \sqrt{f_1 g_1} dM_1 \int_{\Omega_2} \sqrt{f_2 g_2} dM_2 \\ &= \eta_H(P_1, Q_1)\eta_H(P_2, Q_2).\end{aligned}$$

Inductively this implies that

$$\eta_H(P^{\otimes n}, Q^{\otimes n}) = \eta_H(P, Q)^n$$

for arbitrary integers $n \geq 1$.

10.2 Asymptotics for Repeated Binary Experiments

Suppose we observe independent random variables X_1, \dots, X_n with unknown distribution $R \in \{P, Q\}$, where P and Q are two different given probability distributions on (Ω, \mathcal{A}) . Then

$$D_T(P^{\otimes n}, Q^{\otimes n}) \geq D_H(P^{\otimes n}, Q^{\otimes n})^2 = 1 - \eta_H(P, Q)^n$$

converges to 1 exponentially fast. That means, there exists a sequence of tests $\varphi_n : \Omega^n \rightarrow \{0, 1\}$ such that

$$\mathbb{E}_P \varphi_n(\mathbf{X}_n) + \mathbb{E}_Q(1 - \varphi_n(\mathbf{X}_n)) \rightarrow 0,$$

where $\mathbf{X}_n := (X_i)_{i=1}^n$. Throughout this section, asymptotic statements are meant as $n \rightarrow \infty$, unless stated otherwise.

More interesting is the situation when P and Q depend on the sample size n . That means, for each sample size $n \geq 1$ we observe $\mathbf{X}_n = (X_{ni})_{i=1}^n$ with independent components X_{n1}, \dots, X_{nn} having unknown distribution $R_n \in \{P_n, Q_n\}$, where P_n and Q_n are different distributions on (Ω, \mathcal{A}) . The question is, under which conditions on $(P_n)_n$ and $(Q_n)_n$ the potential distributions $P_n^{\otimes n}$ and $Q_n^{\otimes n}$ of \mathbf{X}_n satisfy one of the following three conditions:

- They are (asymptotically) indistinguishable, i.e.

$$D_T(P_n^{\otimes n}, Q_n^{\otimes n}) \rightarrow 0.$$

- They are (asymptotically) perfectly distinguishable, i.e.

$$D_T(P_n^{\otimes n}, Q_n^{\otimes n}) \rightarrow 1.$$

- They are (asymptotically) “interesting” in the sense that

$$\liminf_{n \rightarrow \infty} D_T(P_n^{\otimes n}, Q_n^{\otimes n}) > 0 \quad \text{and} \quad \limsup_{n \rightarrow \infty} D_T(P_n^{\otimes n}, Q_n^{\otimes n}) < 1.$$

It follows from Lemma 10.3 and Exercise 10.8 (a) that the previous three scenarios are equivalent to the analogous ones with D_H or D_H^2 in place of D_T . But note that

$$\begin{aligned} D_H(P_n^{\otimes n}, Q_n^{\otimes n})^2 &= 1 - \eta_H(P_n^{\otimes n}, Q_n^{\otimes n}) \\ &= 1 - \eta_H(P_n, Q_n)^n \\ &= 1 - (1 - D_H(P_n, Q_n)^2)^n, \end{aligned}$$

and the subsequent Lemma 10.11 shows that

$$D_H(P_n^{\otimes n}, Q_n^{\otimes n})^2 = 1 - \exp(-nD_H(P_n, Q_n)^2) + O(n^{-1}).$$

This implies the following results:

Lemma 10.10. (a) For $a \in \{0, 1\}$, the following two conditions are equivalent:

$$\begin{aligned} \lim_{n \rightarrow \infty} D_T(P_n^{\otimes n}, Q_n^{\otimes n}) &= a, \\ \lim_{n \rightarrow \infty} nD_H(P_n, Q_n)^2 &= \begin{cases} \infty & \text{if } a = 1, \\ 0 & \text{if } a = 0. \end{cases} \end{aligned}$$

(b) As $n \rightarrow \infty$, the distance $D_T(P_n^{\otimes n}, Q_n^{\otimes n})$ stays bounded away from 0 and 1 if and only if $nD_H(P_n, Q_n)^2$ stays bounded away from 0 and ∞ .

Lemma 10.11 (Ailam 1968). For arbitrary real numbers $x \in [0, 1]$ and $n \geq 1$,

$$0 \leq e^{-nx} - (1-x)^n \leq e^{-1}n^{-1}.$$

Proof of Lemma 10.11. Since $e^y \geq 1 + y$ for arbitrary $y \in \mathbb{R}$, we know that

$$H_n(x) := e^{-nx} - (1-x)^n$$

satisfies the inequality

$$H_n(x) = (e^{-x})^n - (1-x)^n \geq 0$$

for arbitrary $x \in [0, 1]$. Note also that $H_n(0) = 0$ and $H_n(1) = e^{-n} > 0$. Moreover,

$$H'_n(x) = n(1-x)^{n-1} - ne^{-nx} \begin{cases} \geq 0 & \text{if } n = 1, \\ = 0 & \text{if } n > 1 \text{ and } x = 0, \\ < 0 & \text{if } n > 1 \text{ and } x = 1. \end{cases}$$

Hence

$$\max_{x \in [0,1]} H_1(x) = H_1(0) \leq e^{-1}.$$

For $n > 1$, any maximizer x_n of H_n over $[0, 1]$ has to satisfy $0 < x_n < 1$ and $H'_n(x_n) = 0$, i.e.

$$(1-x_n)^{n-1} = e^{-nx_n}.$$

Consequently,

$$\begin{aligned} \max_{x \in [0,1]} H_n(x) &= H_n(x_n) = e^{-nx_n} - (1-x_n)e^{-nx_n} \\ &= n^{-1} nx_n e^{-nx_n} \\ &\leq n^{-1} \max_{s \geq 0} se^{-s} \\ &= n^{-1} e^{-1}, \end{aligned}$$

because elementary calculations show that se^{-s} is maximized for $s = 1$. □

Expansions of root-densities and log-likelihood ratios. Suppose that for some measure M on (Ω, \mathcal{A}) and arbitrary $n \geq 1$, the densities

$$f_n := \frac{dP_n}{dM} \quad \text{and} \quad g_n := \frac{dQ_n}{dM}$$

exist. Suppose that these densities satisfy the following condition:

(C1) For some probability measure P on (Ω, \mathcal{A}) with density $f = dP/dM$,

$$f_n \rightarrow f \quad \text{in } L^1(M).$$

Furthermore, for some function $h \in L^2(M)$ with $\|h\|_{M,2} > 0$,

$$h_n := \sqrt{n}(\sqrt{g_n} - \sqrt{f_n}) \rightarrow h \quad \text{in } L^2(M).$$

Here and throughout the sequel, $L^r(M)$ is the space of all (equivalence classes of) measurable functions $h : \Omega \rightarrow \mathbb{R}$ such that $\|h\|_{M,r} < \infty$, where

$$\|h\|_{M,r} := \left(\int |h|^r dM \right)^{1/r}$$

for $r \geq 1$. (Two functions h, \tilde{h} are viewed as equivalent if $M(h \neq \tilde{h}) = 0$.)

Some first consequences of this condition:

Lemma 10.12. Under Condition (C1),

$$nD_H(P_n, Q_n)^2 \rightarrow \|h\|_{M,2}^2/2,$$

and

$$D_T(P, P_n), D_T(P, Q_n), D_T(P_n, Q_n) \rightarrow 0.$$

Proof of Lemma 10.12. The convergence of h_n to h in $L^2(M)$ implies that

$$h_n^2 \rightarrow h^2 \quad \text{in } L^1(M),$$

because

$$\begin{aligned} \int |h_n^2 - h^2| dM &= \int |h_n - h| |h_n + h| dM \\ &\leq \|h_n - h\|_{M,2} \|h_n + h\|_{M,2} && \text{(Cauchy-Schwarz)} \\ &\leq \|h_n - h\|_{M,2} (2\|h\|_{M,2} + \|h_n - h\|_{M,2}) && \text{(triangle inequality)} \\ &\rightarrow 0. \end{aligned}$$

In particular,

$$nD_H(P_n, Q_n)^2 = \|h_n\|_{M,2}^2/2 \rightarrow \|h\|_{M,2}^2/2.$$

Note that $f_n \rightarrow f$ in $L^1(M)$ is equivalent to $D_T(P, P_n) \rightarrow 0$, and by Lemma 10.3, $D_T(P_n, Q_n) \leq \sqrt{2} D_H(P_n, Q_n) \rightarrow 0$. Hence by the triangle inequality, $D_T(P, Q_n) \rightarrow 0$ as well. \square

For the next result we have to augment Condition (C1) by an additional one:

(C2) The functions f and h in (C1) satisfy

$$M(h \neq 0 = f) = 0.$$

Lemma 10.13. Suppose that Conditions (C1-2) are satisfied. Let $(A_n)_n$ be an arbitrary sequence of events $A_n \in \mathcal{A}$ such that

$$\min\{P_n(A_n), Q_n(A_n)\} = O(n^{-1}).$$

Then

$$n|Q_n(A_n) - P_n(A_n)| \rightarrow 0.$$

Proof of Lemma 10.13. In terms of the densities f_n and g_n we may write

$$\begin{aligned} n(Q_n(A_n) - P_n(A_n)) &= n \int_{A_n} (g_n - f_n) dM \\ &= n \int_{A_n} (\sqrt{g_n} - \sqrt{f_n})(\sqrt{g_n} + \sqrt{f_n}) dM \\ &= \begin{cases} \int_{A_n} h_n(2\sqrt{nf_n} + h_n) dM, \\ \int_{A_n} h_n(2\sqrt{ng_n} - h_n) dM. \end{cases} \end{aligned}$$

Consequently, by the Cauchy–Schwarz inequality,

$$n|Q_n(A_n) - P_n(A_n)| \leq 2\sqrt{n \min\{P_n(A_n), Q_n(A_n)\} \int_{A_n} h_n^2 dM} + \int_{A_n} h_n^2 dM.$$

Hence it suffices to show that

$$\int_{A_n} h_n^2 dM \rightarrow 0.$$

Since $h_n^2 \rightarrow h^2$ in $L^1(M)$, this is equivalent to

$$\int_{A_n} h^2 dM \rightarrow 0.$$

But it follows from Lemma 10.12 and the assumption that $\min\{P_n(A_n), Q_n(A_n)\} = O(n^{-1})$ that $P(A_n) \rightarrow 0$. Consequently, for any fixed $C > 0$,

$$\int_{A_n} h^2 dM \leq \int_{\{h^2 > Cf\}} h^2 dM + CP(A_n) \rightarrow \int_{\{h^2 > Cf\}} h^2 dM,$$

and

$$\lim_{C \rightarrow \infty} \int_{\{h^2 > Cf\}} h^2 dM = \int_{\{h^2 > 0=f\}} h^2 dM = 0$$

by dominated convergence and our assumption (C2) on f and h . \square

Implications for log-likelihood ratios. Now we consider again the random observation tuple $\mathbf{X}_n = (X_{ni})_{i=1}^n$ with independent components X_{ni} having distribution $R_n \in \{P_n, Q_n\}$. Optimal tests of “ $R_n = P_n$ ” versus “ $R_n = Q_n$ ” are based on the log-likelihood ratio

$$\Lambda_n := \log\left(\frac{dQ_n^{\otimes n}}{dM^{\otimes n}}(\mathbf{X}_n)\right) / \frac{dP_n^{\otimes n}}{dM^{\otimes n}}(\mathbf{X}_n) = \sum_{i=1}^n \log \frac{g_n}{f_n}(X_{ni}) \in [-\infty, \infty]$$

with the conventions that $\log(0) := -\infty$, $\log(\infty) := \infty$, $a/0 := \infty$ for $a > 0$ and $0/0 := 0$. Indeed, since $P_n(f_n > 0) = 1 = Q_n(g_n > 0)$, the random variable Λ_n is well-defined almost surely, and

$$\mathbb{P}_{P_n}(\Lambda_n < \infty) = 1, \quad \mathbb{P}_{Q_n}(\Lambda_n > -\infty) = 1.$$

It may happen with strictly positive probability that $\Lambda_n \in \{-\infty, \infty\}$, but this probability converges to 0. Here is a precise statement:

Theorem 10.14. *Under condition (C),*

$$\Lambda_n \rightarrow_{\mathcal{L}} \begin{cases} \mathcal{N}(-2\|h\|_{M,2}^2, 4\|h\|_{M,2}^2) & \text{if } R_n = P_n \text{ for all } n, \\ \mathcal{N}(+2\|h\|_{M,2}^2, 4\|h\|_{M,2}^2) & \text{if } R_n = Q_n \text{ for all } n. \end{cases}$$

In this result “ $\rightarrow_{\mathcal{L}}$ ” stands for convergence in distribution, meaning that for any fixed continuous function $J : [-\infty, \infty] \rightarrow \mathbb{R}$,

$$\mathbb{E} J(\Lambda_n) \rightarrow \begin{cases} \mathbb{E} J(2\|h\|_{M,2}Z - 2\|h\|_{M,2}^2) & \text{if } R_n = P_n \text{ for all } n \\ \mathbb{E} J(2\|h\|_{M,2}Z + 2\|h\|_{M,2}^2) & \text{if } R_n = Q_n \text{ for all } n \end{cases}$$

with a random variable $Z \sim \mathcal{N}(0, 1)$. Since the limiting distributions are continuous, and since

$$\mathbb{P}(2\|h\|_{M,2}Z \mp 2\|h\|_{M,2}^2 \leq k) = \Phi\left(\frac{k \pm 2\|h\|_{M,2}^2}{2\|h\|_{M,2}}\right),$$

Theorem 10.14 may be rephrased as follows: For arbitrary $k \in \mathbb{R}$,

$$\mathbb{P}(\Lambda_n \leq k) \rightarrow \begin{cases} \Phi\left(\frac{k + 2\|h\|_{M,2}^2}{2\|h\|_{M,2}}\right) & \text{if } R_n = P_n \text{ for all } n, \\ \Phi\left(\frac{k - 2\|h\|_{M,2}^2}{2\|h\|_{M,2}}\right) & \text{if } R_n = Q_n \text{ for all } n. \end{cases}$$

This implies the following result about optimal tests of “ $R_n = P_n$ ” versus “ $R_n = Q_n$ ”:

Corollary 10.15. *Let $\varphi_{n,\alpha} : \Omega^n \rightarrow [0, 1]$ be an optimal level- α test of “ $R_n = P_n$ ” versus “ $R_n = Q_n$ ”. Then*

$$\mathbb{E}_{Q_n} \varphi_{n,\alpha}(\mathbf{X}_n) \rightarrow \Phi(\Phi^{-1}(\alpha) + 2\|h\|_{M,2}).$$

Remark 10.16. Theorem 10.14 and Corollary 10.15 show that under Conditions (C1-2), testing $P_n^{\otimes n}$ versus $Q_n^{\otimes n}$ is asymptotically as difficult as testing

$$\mathcal{N}(0, 1) \text{ versus } \mathcal{N}(\mu, 1),$$

where

$$\mu := 2\|h\|_{M,2}.$$

Indeed,

$$D_H(P_n^{\otimes n}, Q_n^{\otimes n})^2 = 1 - \exp(-nD_H(P_n, Q_n)^2) + O(n^{-1})$$

converges to

$$1 - \exp(-\|h\|_{M,2}^2/2) = 1 - \exp(-\mu^2/8) = D_H(\mathcal{N}(0, 1), \mathcal{N}(\mu, 1))^2.$$

Proof of Corollary 10.15. According to the Neyman–Pearson lemma, we may assume that for some constant $k_{n,\alpha} \in [-\infty, \infty)$,

$$\varphi_{n,\alpha}(\mathbf{X}_n) = \begin{cases} 0 & \text{if } \Lambda_n < k_{n,\alpha}, \\ 1 & \text{if } \Lambda_n > k_{n,\alpha}. \end{cases}$$

But for fixed $k \in \mathbb{R}$,

$$\mathbb{P}_{P_n}(\Lambda_n > k) \rightarrow 1 - \Phi\left(\frac{k + 2\|h\|_{M,2}^2}{2\|h\|_{M,2}}\right).$$

The right hand side is strictly decreasing in k and equals α if, and only if, k is equal to

$$k_\alpha := -2\|h\|_{M,2}^2 + 2\|h\|_{M,2}\Phi^{-1}(1 - \alpha).$$

Hence $k_{n,\alpha} \rightarrow k_\alpha$, and $\mathbb{E}_{Q_n} \varphi_{n,\alpha}(\mathbf{X}_n)$ converges to

$$1 - \Phi\left(\frac{k_\alpha - 2\|h\|_{M,2}^2}{2\|h\|_{M,2}}\right) = 1 - \Phi(-2\|h\|_{M,2} + \Phi^{-1}(1 - \alpha)) = \Phi(\Phi^{-1}(\alpha) + 2\|h\|_{M,2}).$$

□

Proof of Theorem 10.14. It suffices to consider the case $(R_n)_n = (P_n)_n$, because interchanging the roles of $(P_n)_n$ and $(Q_n)_n$ would result in replacing Λ_n with $-\Lambda_n$, and Condition (C) would still be satisfied with $-h$ in place of h .

Since $\sqrt{g_n} = \sqrt{f_n} + h_n/\sqrt{n}$, we may write

$$\Lambda_n = 2 \sum_{i=1}^n \log \frac{\sqrt{g_n}}{\sqrt{f_n}}(X_{ni}) = 2 \sum_{i=1}^n \log \left(1 + \frac{\tilde{h}_n(X_{ni})}{\sqrt{n}} \right)$$

with $\tilde{h}_n := h_n/\sqrt{f_n} \in [-\sqrt{n}, \infty]$. It follows from the wellknown Taylor series of $\log(1+y)$ for $y \in (-1, 1)$ that

$$\log(1+y) = y - \frac{y^2}{2} + \text{rem}(y) \quad \text{with} \quad |\text{rem}(y)| \leq \frac{|y|^3}{3(1-|y|)^+}$$

for arbitrary $y \in [-1, \infty)$. Consequently, with

$$D_n := \max_{1 \leq i \leq n} \frac{|\tilde{h}_n(X_{ni})|}{\sqrt{n}}$$

we obtain the expansion

$$(10.1) \quad \Lambda_n = \frac{2}{\sqrt{n}} \sum_{i=1}^n \tilde{h}_n(X_{ni}) - \frac{1}{n} \sum_{i=1}^n \tilde{h}_n(X_{ni})^2 + \text{Rem}_n$$

with

$$|\text{Rem}_n| \leq \frac{2D_n}{3(1-D_n)^+} \cdot \frac{1}{n} \sum_{i=1}^n \tilde{h}_n(X_{ni})^2.$$

Now we apply the Central Limit Theorem as formulated in Corollary A.24: Suppose we can show that

$$(10.2) \quad \sqrt{n} \int \tilde{h}_n dP_n \rightarrow \mu,$$

$$(10.3) \quad \int \tilde{h}_n^2 dP_n \rightarrow \sigma^2,$$

$$(10.4) \quad \int \tilde{h}_n^2 1_{[\tilde{h}_n^2 \geq n\epsilon]} dP_n \rightarrow 0 \quad \text{for any fixed } \epsilon > 0.$$

Then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{h}_n(X_{ni}) \rightarrow_{\mathcal{L}} \mathcal{N}(\mu, \sigma^2), \quad \frac{1}{n} \sum_{i=1}^n \tilde{h}_n(X_{ni})^2 \rightarrow_p \sigma^2, \quad D_n \rightarrow_p 0.$$

These facts and (10.1) imply that

$$\Lambda_n \rightarrow_{\mathcal{L}} \mathcal{N}(2\mu - \sigma^2, 4\sigma^2).$$

Consequently, it suffices to verify (10.2) with $\mu = -\|h\|_{M,2}^2/2$, (10.3) with $\sigma^2 = \|h\|_{M,2}^2$ and (10.4).

As to (10.2),

$$\begin{aligned}
\sqrt{n} \int \tilde{h}_n dP_n &= n \int \frac{\sqrt{g_n} - \sqrt{f_n}}{\sqrt{f_n}} f_n dM \\
&= n \int (\sqrt{f_n g_n} - f_n) dM \\
&= n(\eta_H(P_n, Q_n) - 1) \\
&= -nD_H(P_n, Q_n)^2 \\
&\rightarrow -\|h\|_{M,2}^2/2,
\end{aligned}$$

see Lemma 10.12. Concerning (10.3),

$$\begin{aligned}
\int \tilde{h}_n^2 dP_n &= \int \frac{h_n^2}{f_n} f_n dM \\
&= \int_{\{f_n > 0\}} h_n^2 dM \\
&= \int_{\{f_n > 0\}} h^2 dM + o(1) \\
&= \int h^2 dM - \int_{\{f_n = 0\}} h^2 dM + o(1) \\
&\rightarrow \|h\|_{M,2}^2,
\end{aligned}$$

because $A_n := \{f_n = 0\}$ satisfies $P(A_n) = P_n(A_n) + o(1) = o(1)$, see also the proof of Lemma 10.13.

It remains to verify (10.4), that means, for any fixed $\epsilon > 0$,

$$\int_{\{h_n^2 \geq \epsilon^2 n f_n\}} h_n^2 dM \rightarrow 0.$$

Again, since $h_n^2 \rightarrow h^2$ in $L^1(M)$, it suffices to show that

$$\int_{\{h_n^2 \geq \epsilon^2 n f_n\}} h^2 dM \rightarrow 0,$$

and the left hand side is equal to

$$\int_0^\infty M(h^2 > r \text{ and } h_n^2 \geq \epsilon^2 n f_n) dr \leq \int_0^\infty M(h^2 > r) dr = \int h^2 dM.$$

Consequently, by dominated convergence, it suffices to show that for any fixed $r > 0$,

$$M(h^2 > r \text{ and } h_n^2 \geq \epsilon^2 n f_n) \rightarrow 0.$$

Indeed, it follows from Markov's inequality that for any fixed $\delta > 0$,

$$\begin{aligned}
&M(h^2 > r \text{ and } h_n^2 \geq \epsilon^2 n f_n) \\
&\leq M(h^2 > r \text{ and } h^2 + \delta \geq \epsilon^2 n(f - \delta)) + M(|h_n^2 - h^2| \geq \delta) + M(|f_n - f| \geq \delta) \\
&\leq M(h^2 > r \text{ and } h^2 + \delta \geq \epsilon^2 n(f - \delta)) + \delta^{-1} \|h_n^2 - h^2\|_{M,1} + \delta^{-1} \|f_n - f\|_{M,1} \\
&\rightarrow M(h^2 > r \text{ and } f \leq \delta).
\end{aligned}$$

Letting $\delta \downarrow 0$, the right hand side converges to $M(h^2 > r \text{ and } f = 0)$, and this equals 0 by assumption (C2). \square

10.3 Fisher Information

Consider a statistical experiment $\mathcal{E} = (\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta})$ with Θ being an open subset of \mathbb{R}^d . Suppose that each P_θ is given by a density $f_\theta > 0$ with respect to some measure M on (Ω, \mathcal{A}) such that for M -almost every $\omega \in \Omega$,

$$\Theta \ni \theta \mapsto f_\theta(\omega)$$

is continuously differentiable with gradient

$$\dot{f}_\theta(\omega) := \left(\frac{\partial f_\theta(\omega)}{\partial \theta_i} \right)_{i=1}^d.$$

Assuming that for any $\theta \in \Theta$, each component of $f_\theta^{-1/2} \dot{f}_\theta$ belongs to $L^2(M)$, the matrix

$$J(\theta) := \int \frac{\dot{f}_\theta \dot{f}_\theta^\top}{f_\theta} dM$$

is well-defined and called the *Fisher information (matrix) of \mathcal{E} at θ* .

Here is yet another interpretation of this matrix: Let $\lambda_\theta := \log f_\theta$. Then for M -almost every $\omega \in \Omega$ the mapping

$$\Theta \ni \theta \mapsto \lambda_\theta(\omega)$$

is continuously differentiable with gradient

$$\dot{\lambda}_\theta(\omega) := \left(\frac{\dot{f}_\theta(\omega)}{f_\theta(\omega)} \right)_{i=1}^d,$$

and

$$J(\theta) = \int \dot{\lambda}_\theta \dot{\lambda}_\theta^\top dP_\theta.$$

Let $(\theta_n)_n$ be a sequence in Θ such that for fixed $\theta \in \Theta$ and $\delta \in \mathbb{R}^d$,

$$\sqrt{n}(\theta_n - \theta) \rightarrow \delta.$$

Then the densities f_{θ_n} and f_θ satisfy

$$\sqrt{n}(\sqrt{f_{\theta_n}} - \sqrt{f_\theta}) \rightarrow h_{\theta, \delta} := \frac{\delta^\top \dot{f}_\theta}{2\sqrt{f_\theta}}$$

almost everywhere. If in addition

$$(10.5) \quad n \int (\sqrt{f_{\theta_n}} - \sqrt{f_\theta})^2 dM \rightarrow \int h_{\theta, \delta}^2 dM,$$

then by Scheffé's theorem, Condition (C) in the previous section is satisfied with P_θ in place of P_n , P_{θ_n} in place of Q_n and limit function $h = h_{\theta, \delta}$. That means, testing $P_\theta^{\otimes n}$ versus $P_{\theta_n}^{\otimes n}$ is asymptotically as difficult as testing $\mathcal{N}(0, 1)$ versus $\mathcal{N}(\mu, 1)$ with

$$\mu = 2\|h\|_{M,2} = \sqrt{\delta^\top J(\theta)\delta}.$$

Condition (10.5) is satisfied, if

$$(10.6) \quad \eta_H(P_\theta, P_{\theta+\delta}) = 1 - \frac{\delta^\top J(\theta)\delta}{8} + o(\|\delta\|^2) \quad \text{as } \delta \rightarrow 0.$$

Definition 10.17 (Regular statistical experiment). A statistical experiment $\mathcal{E} = (\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta})$ satisfying the conditions above is called *regular statistical experiment*.

Example 10.18 (Fisher information in natural exponential families). Suppose that

$$f_\theta = \exp(\theta^\top T - \kappa(\theta))$$

for some measurable mapping $T : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}^d$. Then the experiment \mathcal{E} is regular: Obviously, $f_\theta > 0$. As shown in the exercises, κ is infinitely often differentiable with

$$\kappa(\theta + \delta) = \kappa(\theta) + \delta^\top \mathbb{E}_\theta(T) + \frac{1}{2} \delta^\top \text{Var}_\theta(T) \delta + o(\|\delta\|^2)$$

as $\delta \rightarrow 0$. In particular $f_\theta(\omega)$ is a smooth function of θ with

$$\dot{f}_\theta = (T - \mathbb{E}_\theta(T))f_\theta, \quad \dot{\lambda}_\theta = T - \mathbb{E}_\theta(T),$$

so

$$J(\theta) = \text{Var}_\theta(T).$$

This matrix is positive definite for any $\theta \in \Theta$, unless M^T is concentrated on some hyperplane in \mathbb{R}^d .

Condition (10.6) is satisfied, because

$$\eta_H(P_\theta, P_{\theta+\delta}) = \exp(\kappa(\theta + \delta/2) - \kappa(\theta)/2 - \kappa(\theta + \delta)/2)$$

and

$$\begin{aligned} \kappa(\theta + \delta/2) - \kappa(\theta)/2 - \kappa(\theta + \delta)/2 &= \frac{\delta^\top J(\theta) \delta}{8} - \frac{\delta^\top J(\theta) \delta}{4} + o(\|\delta\|^2) \\ &= -\frac{\delta^\top J(\theta) \delta}{8} + o(\|\delta\|^2). \end{aligned}$$

Remark 10.19 (Smooth transformations of parameters). Let $\mathcal{E} = (\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta})$ be a regular statistical experiment with Fisher information $J(\cdot)$, and let $\tau : \Psi \rightarrow \Omega$ be a diffeomorphism from another open set $\Psi \subset \mathbb{R}^d$ onto Θ . That means, τ is bijective and continuously differentiable with nonsingular Jacobian matrix

$$D\tau(\psi) = \left(\frac{\tau_i(\psi)}{\partial \psi_j} \right)_{i,j=1}^d$$

for any $\psi \in \Psi$.

Then the experiment $\tilde{\mathcal{E}} := (\Omega, \mathcal{A}, (\tilde{P}_\psi)_{\psi \in \Psi})$ with $\tilde{P}_\psi := P_{\tau(\psi)}$ is regular, too, and its Fisher information $\tilde{J}(\cdot)$ is given by

$$\tilde{J}(\psi) = D\tau(\psi)^\top J(\tau(\psi)) D\tau(\psi).$$

This follows from the fact that \tilde{P}_ψ has density $f_{\tau(\psi)}$ with respect to M , and with $\theta = \tau(\psi)$ the chain rule implies that

$$\frac{\partial \tilde{f}_\psi(\omega)}{\partial \psi_j} = \sum_{i=1}^d \frac{\partial f_\theta(\omega)}{\partial \theta_i} \frac{\partial \tau_i(\psi)}{\partial \psi_j} = (D\tau(\psi)^\top \dot{f}_\theta)_j.$$

Moreover, as $\delta \rightarrow 0$,

$$\Delta := \tau(\psi + \delta) - \tau(\psi) = D\tau(\psi)\delta + o(\|\delta\|) = O(\|\delta\|),$$

whence

$$\begin{aligned} \eta_H(\tilde{P}_{\psi+\delta}, \tilde{P}_\psi) &= \eta_H(P_{\theta+\Delta}, P_\theta) \\ &= 1 - \frac{\Delta^\top J(\theta)\Delta}{8} + o(\|\Delta\|^2) \\ &= 1 - \frac{\delta^\top D\tau(\psi)^\top J(\theta)D\tau(\psi)\delta}{8} + o(\|\delta\|^2). \end{aligned}$$

Example 10.20 (Binomial distributions). We observe $X \sim \text{Bin}(n, p)$ with an unknown parameter $p \in (0, 1)$. The natural parameter for the experiment $(\text{Bin}(n, p))_{p \in (0, 1)}$ is given by $\tau(p) := \log(p/(1-p))$ with sufficient statistic $T = X$, and $\tau : (0, 1) \rightarrow \mathbb{R}$ is a diffeomorphism with $\tau'(p) = (p(1-p))^{-1}$. Since $\text{Var}_p(X) = np(1-p)$, Fisher information at p is given by $\tilde{J}(p) = \tau'(p)^2 \text{Var}_p(X)$, i.e.

$$\tilde{J}(p) = \frac{n}{p(1-p)}.$$

Note that $\tilde{J}(p) = \text{Var}_p(\hat{p})^{-1}$ with $\hat{p} := \bar{X}$, which is not a coincidence as explained later.

Example 10.21 (Poisson distributions). We observe $X \sim \text{Poiss}(\lambda)$ with an unknown parameter $\lambda > 0$. The natural parameter for the experiment $(\text{Poiss}(\lambda))_{\lambda > 0}$ is given by $\tau(\lambda) := \log \lambda$ with sufficient statistic $T = X$, and $\tau : (0, \infty) \rightarrow \mathbb{R}$ is a diffeomorphism with $\tau'(\lambda) = \lambda^{-1}$. Since $\text{Var}_\lambda(X) = \lambda$, Fisher information at λ is given by $\tilde{J}(\lambda) = \tau'(\lambda)^2 \text{Var}_\lambda(X)$, i.e.

$$\tilde{J}(\lambda) = \lambda^{-1}.$$

Again $\tilde{J}(\lambda) = \text{Var}_\lambda(\hat{\lambda})^{-1}$ with $\hat{\lambda} := X$.

Implications for point estimation

With our results about testing in Section 10.2 one can prove various precision bounds for point estimators. We present one particular result:

Theorem 10.22 (Asymptotic version of the Cramér–Rao bound). *Let $\mathcal{E} = (\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta})$ be a regular statistical experiment with Fisher information $J(\cdot)$. For each sample size $n \geq 1$ let $\hat{\theta}_n : \Omega^n \rightarrow \mathbb{R}^d$ be an estimator such that for a fixed $\theta \in \Theta$,*

$$\sqrt{n}(\hat{\theta}_n(\mathbf{X}_n) - \theta_n) \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \Sigma(\theta))$$

whenever $\theta_n = \theta + O(n^{-1/2})$ and $\mathbf{X}_n \sim P_{\theta_n}^{\otimes n}$. If $J(\theta)$ is positive definite, then

$$\Sigma(\theta) \geq J(\theta)^{-1}$$

in the sense that $\eta^\top \Sigma(\theta)\eta \geq \eta^\top J(\theta)^{-1}\eta$ for arbitrary $\eta \in \mathbb{R}^d$.

Proof of Theorem 10.22. For a fixed vector $\delta \in \mathbb{R}^d \setminus \{0\}$ define

$$\theta_n := \theta + n^{-1/2}\delta.$$

Then for any vector $\eta \in \mathbb{R}^d \setminus \{0\}$,

$$\sqrt{n}(\widehat{\theta}_n(\mathbf{X}_n) - \theta)^\top \eta \rightarrow_{\mathcal{L}} \begin{cases} \mathcal{N}(0, \eta^\top \Sigma(\theta) \eta) & \text{if } \mathbf{X}_n \sim P_\theta^{\otimes n}, \\ \mathcal{N}(\delta^\top \eta, \eta^\top \Sigma(\theta) \eta) & \text{if } \mathbf{X}_n \sim P_{\theta_n}^{\otimes n}, \end{cases}$$

because

$$\sqrt{n}(\widehat{\theta}_n - \theta)^\top \eta = \sqrt{n}(\widehat{\theta}_n - \theta_n)^\top \eta + \delta^\top \eta.$$

Consequently,

$$\varphi_n := \begin{cases} 1 \\ 0 \end{cases} \quad \text{if } \frac{\sqrt{n}(\widehat{\theta}_n - \theta)^\top \eta}{\sqrt{\eta^\top \Sigma(\theta) \eta}} \begin{cases} \geq \\ \leq \end{cases} \Phi^{-1}(1 - \alpha)$$

defines a test of $P_\theta^{\otimes n}$ versus $P_{\theta_n}^{\otimes n}$ such that

$$\mathbb{E} \varphi_n(\mathbf{X}_n) \rightarrow \begin{cases} \alpha & \text{if } \mathbf{X}_n \sim P_\theta^{\otimes n}. \\ \Phi\left(\Phi^{-1}(\alpha) + \frac{\delta^\top \eta}{\sqrt{\eta^\top \Sigma(\theta) \eta}}\right) & \text{if } \mathbf{X}_n \sim P_{\theta_n}^{\otimes n}. \end{cases}$$

If we would replace φ_n with the optimal level- α test of $P_\theta^{\otimes n}$ versus $P_{\theta_n}^{\otimes n}$, then the asymptotic power under the alternative hypothesis would be

$$\Phi\left(\Phi^{-1}(\alpha) + \sqrt{\delta^\top J(\theta) \delta}\right).$$

Consequently,

$$\frac{\delta^\top \eta}{\sqrt{\eta^\top \Sigma(\theta) \eta}} \leq \sqrt{\delta^\top J(\theta) \delta}.$$

In other words, for arbitrary $\delta, \eta \in \mathbb{R}^d \setminus \{0\}$,

$$\delta^\top \eta \leq \sqrt{\eta^\top \Sigma(\theta) \eta} \sqrt{\delta^\top J(\theta) \delta}.$$

Setting $\delta = J(\theta)^{-1} \eta$ yields the inequality

$$\eta^\top J(\theta)^{-1} \eta \leq \eta^\top \Sigma(\theta) \eta$$

for arbitrary $\eta \in \mathbb{R}^d \setminus \{0\}$. □

Example 10.23 (Maximum-likelihood estimation in natural exponential families). Let \mathcal{E} be a natural exponential family with sufficient statistic $T : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}^d$ such that M^T is not concentrated on a hyperplane in \mathbb{R}^d . Let $(\theta_n)_n$ be a sequence in Θ with limit $\theta \in \Theta$, and let $\mathbf{X}_n \sim P_{\theta_n}^{\otimes n}$. Then the log-likelihood function

$$L_n = L_n(\cdot, \mathbf{X}_n) : \Theta \rightarrow \mathbb{R}, \quad L_n(\theta) := \sum_{i=1}^n \log f_\theta(X_{ni})$$

has the following property: With probability tending to 1, there exists a unique maximizer $\hat{\theta}_n = \hat{\theta}_n(\mathbf{X}_n)$ of L_n , and

$$\sqrt{n}(\hat{\theta}_n - \theta_n) \rightarrow \mathcal{N}_d(0, J(\theta)^{-1}) = \mathcal{N}_d(0, \text{Var}_\theta(T)^{-1}).$$

To verify these claims, note first that

$$L_n(\theta) = n(\bar{T}_n^\top \theta - \kappa(\theta))$$

with $\bar{T}_n := n^{-1} \sum_{i=1}^n T(X_{ni})$, and thus

$$\begin{aligned} \nabla L_n(\theta) &= n(\bar{T}_n - \mathbb{E}_\theta(T)), \\ D^2 L_n(\theta) &= -n \text{Var}_\theta(T). \end{aligned}$$

This shows that L_n is strictly concave, whence L_n has a unique maximizer or no maximizer at all.

One can deduce from the multivariate version of Lindeberg's Central Limit Theorem that

$$Z_n := \sqrt{n}(\bar{T}_n - \mathbb{E}_{\theta_n}(T)) \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \text{Var}_{\theta_n}(T)).$$

Now we introduce the localized log-likelihood function $H_n : \mathbb{R}^d \rightarrow [-\infty, \infty)$ with

$$H_n(\delta) := L_n(\theta_n + n^{-1/2}\delta) - L_n(\theta_n),$$

where $H_n(\delta) = -\infty$ if, and only if, $\kappa(\theta_n + n^{-1/2}\delta) = \infty$. Then elementary calculations reveal that

$$H_n(\delta) = Z_n^\top \delta - \frac{1}{2} \delta^\top \text{Var}_\theta(T) \delta + \text{Rem}_n(\delta)$$

where

$$\sup_{\delta: \|\delta\| \leq C} |\text{Rem}_n(\delta)| \rightarrow_p 0$$

for any fixed $C > 0$. From this one may deduce that with probability tending to one, H_n has a unique maximizer given by

$$\text{Var}_\theta(T)^{-1} Z_n + o_p(1).$$

But this is equivalent to saying that with asymptotic probability 1, the unique maximizer $\hat{\theta}_n$ of L_n exists and satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_n) = \text{Var}_\theta(T)^{-1} Z_n + o_p(1).$$

In particular,

$$\sqrt{n}(\hat{\theta}_n - \theta_n) \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \text{Var}_\theta(T)^{-1})$$

because $Z_n \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \text{Var}_\theta(T))$.

Example 10.24 (Maximum-likelihood estimation in smoothly parametrized exponential families). Let $\tilde{\mathcal{E}} = (\Omega, \mathcal{A}, (\tilde{P}_\psi)_{\psi \in \Psi})$ be an exponential family with sufficient statistic $T : \Omega \rightarrow \mathbb{R}^d$, i.e. \tilde{P}_ψ has density

$$\tilde{f}_\psi = \exp(\tau(\psi)^\top T - \kappa(\tau(\psi)))$$

for some bijective mapping $\tau : \Psi \rightarrow \Theta$, where Ψ and Θ are open subsets of \mathbb{R}^d , and

$$\kappa(\theta) := \log \int \exp(\theta^\top T) dM.$$

Suppose further that τ is a diffeomorphism.

Now let $(\psi_n)_n$ be a sequence in Ψ with limit $\psi \in \Psi$, and let $\mathbf{X}_n \sim \tilde{P}_{\psi_n}^{\otimes n}$. Then the log-likelihood function $L_n = L_n(\cdot, \mathbf{X}_n) := \sum_{i=1}^n \log \tilde{f}_{\psi}(X_{ni})$ has the following property: With asymptotic probability 1, there exists a unique maximizer $\hat{\psi}_n = \hat{\psi}_n(\mathbf{X}_n)$ of $L_n(\cdot)$, and

$$\sqrt{n}(\hat{\psi}_n - \psi_n) \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \tilde{J}(\psi)^{-1})$$

with the Fisher information $\tilde{J}(\cdot)$ of $\tilde{\mathcal{E}}$, i.e. $\tilde{J}(\psi) = D\tau(\psi)^\top \text{Var}_{\psi}(T) D\tau(\psi)$.

With $\theta_n := \tau(\psi_n)$ and $\theta := \tau(\psi)$, it follows from Example 10.23 that with asymptotic probability 1 there exists a unique maximum likelihood estimator $\hat{\theta}_n = \hat{\theta}_n(\mathbf{X}_n)$ for the experiment $\mathcal{E} = (\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta})$ with $P_{\tau(\psi)} := \tilde{P}_\psi$ such that

$$\sqrt{n}(\hat{\theta}_n - \theta_n) \rightarrow_{\mathcal{L}} \mathcal{N}_d(0, \text{Var}_{\psi}(T)^{-1}).$$

But then $\hat{\psi}_n := \tau^{-1}(\hat{\theta}_n)$ is a maximum likelihood estimator for $\tilde{\mathcal{E}}$, and elementary calculus reveals that

$$\begin{aligned} \sqrt{n}(\hat{\psi}_n - \psi_n) &= D\tau(\psi)^{-1} \sqrt{n}(\hat{\theta}_n - \theta_n) + o_p(1) \\ &\rightarrow_{\mathcal{L}} \mathcal{N}_d(0, D\tau(\psi)^{-1} \text{Var}_{\psi}(T) (D\tau(\psi)^{-1})^\top) \\ &= \mathcal{N}_d(0, \tilde{J}(\psi)^{-1}). \end{aligned}$$

Appendix A

Auxiliary Results

A.1 Some Basics from Measure Theory

For reference in the main text, we collect some basic notions and results in this section; for a thorough account of measure theory, we refer to the monograph of Bauer (2001). Throughout this section let Ω be a nonvoid set, and let $\mathcal{P}(\Omega)$ be the family of all subsets of Ω .

Fields, σ -fields, contents and measures

Definition A.1 (Field and σ -field). A family $\mathcal{A} \subset \mathcal{P}(\Omega)$ is called a field over Ω , if it satisfies the following three conditions:

(F.1) $\emptyset, \Omega \in \mathcal{A}$.

(F.2) If $A \in \mathcal{A}$, then $\Omega \setminus A \in \mathcal{A}$.

(F.3) If $A, B \in \mathcal{A}$, then $A \cup B, A \cap B \in \mathcal{A}$.

A family $\mathcal{A} \subset \mathcal{P}(\Omega)$ is called a σ -field over Ω , if it satisfies (F.1-2) and the following condition:

(F.3) $_{\sigma}$ If $A_1, A_2, A_3, \dots \in \mathcal{A}$, then $\bigcup_{n \geq 1} A_n, \bigcap_{n \geq 1} A_n \in \mathcal{A}$.

Remark A.2. If condition (F.2) is satisfied, then (F.1) is equivalent to $\Omega \in \mathcal{A}$, and in conditions (F.3) and (F.3) $_{\sigma}$ it suffices to consider only the unions or only the intersections of sets in \mathcal{A} .

Definition A.3 (Content and measure). Let \mathcal{A} be a field over Ω . A mapping $M : \mathcal{A} \rightarrow [0, \infty]$ is called a *content*, if it satisfies the following two conditions:

(M.1) $M(\emptyset) = 0$.

(M.2) $M(A \cup B) = M(A) + M(B)$ for disjoint sets $A, B \in \mathcal{A}$.

The mapping M is called a *measure*, if it satisfies (M.1) and the following condition:

(M.2) $_{\sigma}$ $M(\bigcup_{n \geq 1} A_n) = \sum_{n \geq 1} M(A_n)$ for pairwise disjoint sets A_1, A_2, A_3, \dots in \mathcal{A} such that $\bigcup_{n \geq 1} A_n \in \mathcal{A}$.

A content or measure M is called *finite* if $M(\Omega) < \infty$. In case of $M(\Omega) = 1$, it is called *probability content* or *probability measure*, respectively.

Any measure is obviously a content. More interesting is the question under which additional property a content is a measure.

Lemma A.4. *A content M on a field \mathcal{A} is a measure if and only if it satisfies the following continuity property: If $B_1 \subset B_2 \subset B_3 \subset \dots$ are sets in \mathcal{A} such that $B := \bigcup_{n \geq 1} B_n \in \mathcal{A}$, then $M(B) = \lim_{n \rightarrow \infty} M(B_n)$.*

A finite content M on a field \mathcal{A} is a measure if and only if $\lim_{n \rightarrow \infty} M(A_n) = 0$ for arbitrary sets $A_1 \supset A_2 \supset A_3 \supset \dots$ in \mathcal{A} with $\bigcap_{n \geq 1} A_n = \emptyset$.

Here is an important result about the extension of measures:

Theorem A.5 (Carathéodory). *Let M be a finite measure on a field \mathcal{A} over Ω , and let \mathcal{A}_* be the smallest σ -field over Ω containing \mathcal{A} . Then there exists a unique measure M_* on \mathcal{A}_* such that $M_*(A) = M(A)$ for all $A \in \mathcal{A}$.*

Measurability and integrals

Measurability. Let (Ω, \mathcal{A}) and $(\mathcal{X}, \mathcal{B})$ be measurable spaces, that means, \mathcal{A} is a σ -field over Ω , and \mathcal{B} is a σ -field over \mathcal{X} . A mapping $X : \Omega \rightarrow \mathcal{X}$ is called \mathcal{A} - \mathcal{B} -measurable if

$$X^{-1}(B) := \{\omega \in \Omega : X(\omega) \in B\} \in \mathcal{A} \quad \text{for all } B \in \mathcal{B}.$$

Sometimes one abuses notation slightly and talks about a measurable function $X : (\Omega, \mathcal{A}) \rightarrow (\mathcal{X}, \mathcal{B})$.

Suppose that \mathcal{B} is generated by some family $\mathcal{E} \subset \mathcal{P}(\mathcal{X})$. That means, $\mathcal{B} = \sigma(\mathcal{E})$ is the smallest σ -field containing \mathcal{E} . Then X is \mathcal{A} - \mathcal{B} -measurable if and only if

$$X^{-1}(E) \in \mathcal{A} \quad \text{for all } E \in \mathcal{E}.$$

If it is clear from the context which σ -field \mathcal{B} the set \mathcal{X} is equipped with, one talks about a \mathcal{A} -measurable mapping $X : \Omega \rightarrow \mathcal{X}$ or a measurable mapping $X : (\Omega, \mathcal{A}) \rightarrow \mathcal{X}$. In particular, if \mathcal{X} is some interval in $\bar{\mathbb{R}} = [-\infty, \infty]$, then \mathcal{B} is tacitly understood to be the Borel- σ -field, i.e. the smallest σ -field containing all intervals in \mathcal{X} .

If $(f_n)_{n \geq 1}$ is a sequence of \mathcal{A} -measurable functions $f_n : \Omega \rightarrow \bar{\mathbb{R}}$ which converges pointwise to a function $f : \Omega \rightarrow \bar{\mathbb{R}}$, then f is \mathcal{A} -measurable as well.

Lebesgue integrals. Let M be a measure on a σ -field \mathcal{A} over Ω . For an \mathcal{A} -measurable function $f : \Omega \rightarrow \bar{\mathbb{R}}$, its integral with respect to M is defined as follows:

Case 1. Let \mathcal{G}_+ be the set of all functions f of the form $f = \sum_{i=1}^n \lambda_i 1_{A_i}$ with $n \in \mathbb{N}$, $\lambda_i \in [0, \infty)$ and $A_i \in \mathcal{A}$. For such a function we define

$$\int f dM := \sum_{i=1}^n \lambda_i M(A_i).$$

One can show that this definition does not depend on the particular representation of f .

Case 2. Let $f : \Omega \rightarrow [0, \infty]$ be \mathcal{A} -measurable. Then

$$\int f dM := \sup \left\{ \int g dM : g \in \mathcal{G}_+, g \leq f \right\}.$$

Two important properties of this integral are: For \mathcal{A} -measurable functions $f, g : \Omega \rightarrow [0, \infty]$,

$$\int (af + bg) dM = a \int f dM + b \int g dM \quad \text{for arbitrary constants } a, b \in [0, \infty],$$

and

$$\int f dM \leq \int g dM \quad \text{if } f \leq g.$$

Case 3. Let $f : \Omega \rightarrow \bar{\mathbb{R}}$ be \mathcal{A} -measurable. We write $f = f^+ - f^-$ with $f^\pm := \max(\pm f, 0)$, so $|f| = f^+ + f^-$. Then

$$\int f dM := \int f^+ dM - \int f^- dM,$$

provided that one of the two integrals $\int f^\pm dM$ is finite. Otherwise the integral $\int f dM$ is not defined. This definition implies that

$$\left| \int f dM \right| \leq \int |f| dM.$$

Moreover, $\int f dM$ is well-defined in \mathbb{R} if and only if $\int |f| dM$ is finite.

Sometimes it is useful to indicate arguments of the functions to be integrated, so we also write

$$\int f dM = \int f(\omega) M(d\omega).$$

In addition, for sets $A \in \mathcal{A}$ one often writes

$$\int_A f dM := \int 1_A f dM,$$

so $\int f dM = \int_\Omega f dM$.

Monotone convergence. If $(f_n)_{n \geq 1}$ is a sequence of \mathcal{A} -measurable functions $f_n : \Omega \rightarrow [0, \infty]$, and is $f_n \uparrow f$ pointwise as $n \rightarrow \infty$, then

$$\lim_{n \rightarrow \infty} \int f_n dM = \int f dM.$$

Dominated convergence. Let $(f_n)_{n \geq 1}$ be a sequence of \mathcal{A} -measurable functions $f_n : \Omega \rightarrow \bar{\mathbb{R}}$ converging pointwise to a function f . Suppose that $|f_n| \leq g$ for some \mathcal{A} -measurable function $g : \Omega \rightarrow [0, \infty]$ such that $\int g dM < \infty$. Then all integrals $\int f_n dM$ and $\int f dM$ are well-defined in \mathbb{R} , and

$$\lim_{n \rightarrow \infty} \int |f_n - f| dM = 0, \quad \lim_{n \rightarrow \infty} \int f_n dM = \int f dM.$$

Scheffé's theorem. Let $(f_n)_{n \geq 1}$ be a sequence of \mathcal{A} -measurable functions $f_n : \Omega \rightarrow \overline{\mathbb{R}}$ converging pointwise to a function f . Suppose that

$$\limsup_{n \rightarrow \infty} \int |f_n| dM \leq \int |f| dM < \infty.$$

Then the integrals $\int f_n dM$ and $\int f dM$ are well-defined in \mathbb{R} for sufficiently large n , and

$$\lim_{n \rightarrow \infty} \int |f_n - f| dM = 0, \quad \lim_{n \rightarrow \infty} \int f_n dM = \int f dM.$$

Dynkin systems

To verify measurability of certain functions or identity of two measures, the following type of set families is very useful.

Definition A.6 (Dynkin system). A family \mathcal{D} of subsets of Ω is called a *Dynkin system* over Ω if the following three conditions are satisfied:

(D.1) $\Omega \in \mathcal{D}$.

(D.2) If $A \in \mathcal{D}$, then $\Omega \setminus A \in \mathcal{D}$.

(D.3) If A_1, A_2, A_3, \dots are pairwise disjoint sets in \mathcal{D} , then $\bigcup_{n \geq 1} A_n \in \mathcal{D}$.

Remark A.7. Conditions (D.1–3) are equivalent to the following three conditions:

(D'.1) $\Omega \in \mathcal{D}$.

(D'.2) If $A, B \in \mathcal{D}$ with $A \subset B$, then $B \setminus A \in \mathcal{D}$;

(D'.3) If $B_1 \subset B_2 \subset B_3 \subset \dots$ are sets in \mathcal{D} , then $\bigcup_{n \geq 1} B_n \in \mathcal{D}$.

Remark A.8. Let \mathcal{D} be a Dynkin system over Ω . Then \mathcal{D} is a σ -field over Ω if and only if it is closed under (finite) intersections, that means, $A \cap B \in \mathcal{D}$ for arbitrary $A, B \in \mathcal{D}$.

Recall also the following well-known fact about Dynkin systems.

Theorem A.9 (Dynkin). *Let \mathcal{E} be an arbitrary family of subsets of Ω . There exists a smallest Dynkin system $\mathcal{D}(\mathcal{E})$ over Ω such that $\mathcal{E} \subset \mathcal{D}(\mathcal{E})$. If \mathcal{E} is closed under (finite) intersections, then the Dynkin system $\mathcal{D}(\mathcal{E})$ coincides with the smallest σ -field $\sigma(\mathcal{E})$ containing \mathcal{E} .*

A classical application of Dynkin systems is the following result about uniqueness of measures.

Theorem A.10 (Dynkin). *Suppose that P and Q are probability measures on (Ω, \mathcal{A}) , where $\mathcal{A} = \sigma(\mathcal{E})$ for some family $\mathcal{E} \subset \mathcal{P}(\Omega)$. Suppose further that $P \equiv Q$ on \mathcal{E} . If \mathcal{E} is closed under (finite) intersections, then $Q \equiv P$ on \mathcal{A} .*

Proof of Theorem A.10. One can easily verify that $\mathcal{D} := \{A \in \mathcal{A} : P(A) = Q(A)\}$ defines a Dynkin system containing \mathcal{E} . Thus $\mathcal{D}(\mathcal{E}) \subset \mathcal{D} \subset \mathcal{A}$. But \cap -stability of \mathcal{E} implies that $\mathcal{D}(\mathcal{E})$ coincides with $\sigma(\mathcal{E}) = \mathcal{A}$, see Theorem A.9. Hence, $\mathcal{D} = \mathcal{A}$. \square

Exercise A.11. Let $\Omega = \{1, 2, 3, 4\}$, and let $\mathcal{E} = \{\{1, 2\}, \{1, 4\}\}$.

(a) Show that $\sigma(\mathcal{E}) = \mathcal{P}(\Omega)$.

- (b) Determine $\mathcal{D}(\mathcal{E})$.
- (c) Find two different probability measures P, Q on $\mathcal{P}(\Omega)$ such that $P \equiv Q$ on $\mathcal{D}(\mathcal{E})$.

A.2 Two Compactness Properties of Statistical Tests

Let M be a σ -finite measure on a measurable space (Ω, \mathcal{A}) . Further, let \mathcal{F} be the set of all \mathcal{A} -measurable functions $f : \Omega \rightarrow \mathbb{R}$ with $\int |f| dM < \infty$, and let \mathcal{T} be the set of all \mathcal{A} -measurable functions $\varphi : \Omega \rightarrow [0, 1]$. Then the set \mathcal{T} satisfies the following compactness condition:

Theorem A.12 (Weak compactness of \mathcal{T}). *The set*

$$\left\{ \left(\int \varphi f dM \right)_{f \in \mathcal{F}} : \varphi \in \mathcal{T} \right\}$$

is a convex and compact subset of $\mathbb{R}^{\mathcal{F}}$, where the latter set is equipped with the usual product topology.

The set $\mathbb{R}^{\mathcal{F}}$ is the set of all tuples $(x_f)_{f \in \mathcal{F}}$ with components $x_f \in \mathbb{R}$. The product topology on this set is the smallest topology such that the mapping $\mathbb{R}^{\mathcal{F}} \ni x \mapsto x_f \in \mathbb{R}$ is continuous for arbitrary $f \in \mathcal{F}$.

Exercise A.13. The proof of Theorem A.12 relies on Tikhonov's theorem, hence on the axiom of choice. Prove a simpler result without this tool in the special case of a countable set Ω and M being the counting measure on Ω : For arbitrary $m \in \mathbb{N}$ and functions $f_1, \dots, f_m \in \mathcal{L}^1(M)$, the set

$$\left\{ \left(\int \varphi f_j dM \right)_{j=1}^m : \varphi \in \mathcal{T} \right\}$$

is a convex and compact subset of \mathbb{R}^m .

Proof of Theorem A.12. Writing $f \in \mathcal{F}$ as $f^+ - f^-$ with $f^\pm := \max(\pm f, 0)$, we know that

$$\int \varphi f dM \in K_f := \left[- \int f^- dM, \int f^+ dM \right]$$

for any $\varphi \in \mathcal{T}$ and $f \in \mathcal{F}$. Hence,

$$\mathcal{K}_* := \left\{ \left(\int \varphi f dM \right)_{f \in \mathcal{F}} : \varphi \in \mathcal{T} \right\}$$

is a subset of

$$\mathcal{K} := \{x \in \mathbb{R}^{\mathcal{F}} : x_f \in K_f \text{ for all } f \in \mathcal{F}\}.$$

Since each K_f , $f \in \mathcal{F}$, is a compact interval, it follows from Tikhonov's theorem that \mathcal{K} is a compact subset of $\mathbb{R}^{\mathcal{F}}$.

Linearity of integrals implies that \mathcal{K}_* is a subset of

$$\mathcal{K}_o := \{x \in \mathcal{K} : f \mapsto x_f \text{ is linear}\}.$$

That means, \mathcal{K}_o consists of all tuples $x \in \mathcal{K}$ such that for arbitrary $f, g \in \mathcal{F}$ and $\lambda \in \mathbb{R}$,

$$\begin{aligned} x_{\lambda f} &= \lambda x_f, \\ x_{f+g} &= x_f + x_g. \end{aligned}$$

Note that each of the previous two constraints defines a closed subset of $\mathbb{R}^{\mathcal{F}}$. Hence the set \mathcal{K}_o is a closed subset of \mathcal{K} , i.e. it is compact. Moreover, one can easily verify that \mathcal{K}_o is convex.

Now the assertion of Theorem A.12 is true if we can show that $\mathcal{K}_* = \mathcal{K}_o$. That means, we have to show that for any fixed $x \in \mathcal{K}_o$ there exists a test $\varphi \in \mathcal{T}$ such that $x_f = \int \varphi f dM$ for arbitrary $f \in \mathcal{F}$. Indeed, it follows from linearity of $f \mapsto x_f$ and the inclusion $x_f \in K_f$ for all $f \in \mathcal{F}$ that $f \mapsto x_f$ defines a continuous linear functional on $\mathcal{L}^1(M)$, equipped with the seminorm $\|f\| := \int |f| dM$. Consequently, by Theorem 2.29, there exists a bounded measurable function $\varphi : \Omega \rightarrow \mathbb{R}$ such that $x_f = \int f \varphi dM$ for all $f \in \mathcal{F}$. In particular, since $0 \leq x_{1_A} = \int_A \varphi dM \leq M(A)$ for all $A \in \mathcal{A}$, the function φ satisfies $M(\{\varphi < 0\} \cup \{\varphi > 1\}) = 0$. Hence, we may assume that $\varphi \in \mathcal{T}$. \square

The next result establishes a sequential compactness property of \mathcal{T} . Its proof is constructive in the sense that it does not use the axiom of choice.

Theorem A.14 (Weak sequential compactness of \mathcal{T}). *Let $(\varphi_n)_{n \geq 1}$ be a sequence in \mathcal{T} . Then there exist a subsequence $(\varphi_{n(k)})_{k \geq 1}$ and a test $\varphi \in \mathcal{T}$ such that*

$$\lim_{k \rightarrow \infty} \int \varphi_{n(k)} f dM = \int \varphi f dM \quad \text{for any } f \in \mathcal{F}.$$

Proof of Theorem A.14. As in the proof of Theorem A.12, one can reduce the claim to the case of a probability measure M . Let \mathcal{A}_o be the σ -field generated by the tests φ_n , $n \geq 1$. This sub- σ -field of \mathcal{A} has a countable generator, for instance, the family of all sets $\{\varphi_n \leq q\}$, $n \in \mathbb{N}$, $q \in \mathbb{Q}$. Consequently, there exists a countable field \mathcal{A}_{oo} of Ω such that $\mathcal{A}_o = \sigma(\mathcal{A}_{oo})$. By Cantor's diagonalisation trick, there exists a subsequence $(\varphi_{n(k)})_{k \geq 1}$ such that

$$L(1_{A_{oo}}) := \lim_{k \rightarrow \infty} \int \varphi_{n(k)} 1_{A_{oo}} dM$$

exists for all $A_{oo} \in \mathcal{A}_{oo}$.

Let \mathcal{F}_* be the set of all $f \in \mathcal{F}$ such that the limit

$$L(f) := \lim_{k \rightarrow \infty} \int \varphi_{n(k)} f dM$$

exists. One can easily verify that \mathcal{F}_* is a linear subspace of \mathcal{F} , and L is continuous on \mathcal{F}_* with respect to the seminorm $f \mapsto \|f\| := \int |f| dM$. Precisely, for any $f \in \mathcal{F}$ and $n \geq 1$,

$$-\int f^- dM \leq \int \varphi_n f dM \leq \int f^+ dM,$$

and thus,

$$(A.1) \quad -\int f^- dM \leq L(f) \leq \int f^+ dM \quad \text{for all } f \in \mathcal{F}_*.$$

The space \mathcal{F}_* is also closed with respect to $\|\cdot\|$. For if $(f_\ell)_{\ell \geq 1}$ is a sequence in \mathcal{F}_* with limit $f \in \mathcal{F}$, then for any fixed $\ell \geq 1$,

$$\begin{aligned} & \limsup_{k, k' \rightarrow \infty} \left| \int \varphi_{n(k)} f \, dM - \int \varphi_{n(k')} f \, dM \right| \\ & \leq 2\|f - f_\ell\| + \limsup_{k, k' \rightarrow \infty} \left| \int \varphi_{n(k)} f_\ell \, dM - \int \varphi_{n(k')} f_\ell \, dM \right| = 2\|f - f_\ell\|, \end{aligned}$$

and the right hand side tends to 0 as $\ell \rightarrow \infty$. Consequently, $(\int \varphi_{n(k)} f \, dM)_{k \geq 1}$ is a Cauchy sequence in \mathbb{R} .

The space \mathcal{F}_* contains all indicator functions $1_{A_{oo}}$, $A_{oo} \in \mathcal{A}_{oo}$, and the linear span of the latter functions is dense in $\mathcal{L}^1(M|_{\mathcal{A}_o})$ with respect to $\|\cdot\|$. Thus, \mathcal{F}_* contains all functions in $\mathcal{L}^1(M|_{\mathcal{A}_o})$. But for an arbitrary function $f \in \mathcal{F}$ and its conditional expectation $f_o := \mathbb{E}(f | \mathcal{A}_o)$ with respect to the probability measure M ,

$$\int \varphi_n f \, dM = \int \varphi_n f_o \, dM$$

for all $n \geq 1$, so $\mathcal{F}_* = \mathcal{F}$.

It follows from (A.1) that L is a continuous linear functional on $\mathcal{L}^1(M)$, so Theorem 2.29 implies the existence of a bounded measurable function $\varphi : \Omega \rightarrow \mathbb{R}$ such that $L(f) = \int f \varphi \, dM$ for all $f \in \mathcal{F}$. In particular, since $0 \leq L(1_A) = \int_A \varphi \, dM \leq M(A)$ for all $A \in \mathcal{A}$, the function φ satisfies $M(\{\varphi < 0\} \cup \{\varphi > 1\}) = 0$. Hence, we may assume that $\varphi \in \mathcal{T}$. \square

A.3 Uniqueness of Moment-Generating Functions

In the context of completeness of statistical experiments and exponential families we utilize a classical result from measure theory.

Theorem A.15. *Let M be a measure on \mathbb{R}^d , and let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a measurable function such that*

$$\int_{\mathbb{R}^d} \exp(u^\top x) f(x) M(dx) = 0$$

for all u in a nonempty open set $U \subset \mathbb{R}^d$. Then

$$M(f \neq 0) = 0.$$

Proof of Theorem A.15. Suppose first that $0 \in \mathbb{R}^d$ is an interior point of U . Then for some $\epsilon > 0$,

$$L(u) := \int \exp(u^\top x) f(x) M(dx) = 0 \quad \text{for all } u \in (-\epsilon, \epsilon)^d.$$

One can easily verify that $L(u)$ is well-defined in \mathbb{C} for all complex vectors

$$u \in U_o^d \quad \text{with} \quad U_o := \{z \in \mathbb{C} : -\epsilon < \operatorname{Re} z < \epsilon\}.$$

Moreover, for any given index $j \in \{1, \dots, d\}$, $L(u)$ is a holomorphic (i.e. complex differentiable) function of $u_j \in U_o$ while $(u_k)_{k \neq j}$ is fixed; see Exercise A.16. But it is well-known from complex

analysis that a holomorphic function $H : U_o \rightarrow \mathbb{C}$ with $H \equiv 0$ on $(-\epsilon, \epsilon)$ satisfies $H \equiv 0$ on U_o . Consequently, we may conclude inductively for $j = 1, 2, \dots, d$ that

$$\begin{aligned} L \equiv 0 \text{ on } (-\epsilon, \epsilon)^d & \quad \text{whence} \quad L \equiv 0 \text{ on } U_o \times (-\epsilon, \epsilon)^{d-1} \\ & \quad \text{whence} \quad L \equiv 0 \text{ on } U_o^2 \times (-\epsilon, \epsilon)^{d-2} \\ & \quad \dots \\ & \quad \text{whence} \quad L \equiv 0 \text{ on } U_o^d. \end{aligned}$$

Since $\{iy : y \in \mathbb{R}^d\} \subset U_o^d$ (with $i = \sqrt{-1}$), the latter equality for L implies that

$$\int \exp(iy^\top x) f^+(x) M(dx) = \int \exp(iy^\top x) f^-(x) M(dx) \quad \text{for all } y \in \mathbb{R}^d.$$

That means, the characteristic functions of the finite measures Q^+ and Q^- , where

$$Q^\pm(A) := \int_A f^\pm dM,$$

are identical. But a finite measure is uniquely determined by its characteristic function, so $Q^+ \equiv Q^-$. In particular, $Q^\pm(\mathbb{R}^d) = Q^\pm(f^\pm > 0) = Q^\mp(f^\pm > 0) = 0$, and this implies that $M(f \neq 0) = 0$.

In case of a general open set U , let u_o be an interior point of U . Then $V := U - u_o$ is an open neighborhood of 0, and with $g(x) := \exp(u_o^\top x) f(x)$ the assumption reads

$$\int \exp(v^\top x) g(x) M(dx) = 0 \quad \text{for all } v \in V.$$

But then the previous considerations show that $M(f \neq 0) = M(g \neq 0) = 0$. □

Exercise A.16. Let (Ω, \mathcal{A}, M) be a measure space, $g : (\Omega, \mathcal{A}) \rightarrow \mathbb{C}$ and $T : (\Omega, \mathcal{A}) \rightarrow \mathbb{R}$ be measurable, and suppose that for real numbers $a < b$,

$$\int |g| \exp(cT) dM < \infty \quad \text{for } c = a, b.$$

Show that

$$f(z) := \int g \exp(zT) dM$$

defines a holomorphic function on $\{z \in \mathbb{C} : a < \operatorname{Re}(z) < b\}$.

A.4 Hoeffding's Decomposition

Hoeffding's decomposition is a generalization of Hájek's projection as described in Lemma 8.38. The setting is the same, we consider a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with stochastically independent random variables X_1, \dots, X_n with values in $(\mathcal{X}_1, \mathcal{B}_1), \dots, (\mathcal{X}_n, \mathcal{B}_n)$, respectively. Now we consider the Hilbert space \mathbb{H} of all random variables $Y \in L^2(\mathbb{P})$ which are a measurable function of the random tuple $X := (X_1, \dots, X_n)$ with values in $\mathcal{X} := \mathcal{X}_1 \times \dots \times \mathcal{X}_n$.

For any nonvoid set $K \subset \{1, \dots, n\}$ let \mathbb{H}_K be the subspace of all random variables $Y \in \mathbb{H}$ which are a measurable function of

$$X_K := (X_i)_{i \in K},$$

and let \mathbb{H}_\emptyset be the subspace of all constant random variables. In particular, $\mathbb{H} = \mathbb{H}_{\{1, \dots, n\}}$. The orthogonal projection of \mathbb{H} onto \mathbb{H}_K is given by Π_K with

$$\Pi_K Y := \begin{cases} \mathbb{E}(Y) & \text{if } K = \emptyset \\ \mathbb{E}(Y | X_K) & \text{else} \end{cases}$$

for $Y \in \mathbb{H}$. Strictly speaking, we should write $\mathbb{E}(Y | \sigma(X_K))$, but $\mathbb{E}(Y | X_K)$ is more convenient and intuitive. One treats X_K temporarily as a fixed tuple, and if $Y = f(X_K, X_L)$ with $L = \{1, \dots, n\} \setminus K$, then

$$\mathbb{E}(Y | X_K) = \int f(X_K, z) P_L(dz),$$

where P_L denotes the distribution of X_L .

A key property of these projections Π_K is that

$$(A.2) \quad \Pi_J \Pi_K = \Pi_{J \cap K} \quad \text{for arbitrary } J, K \subset \{1, \dots, n\}.$$

This can be easily derived from Fubini's theorem. In particular, $\Pi_J \Pi_K = \Pi_K \Pi_J$. Now we define

$$\Pi_K^o := \sum_{I \subset K} (-1)^{\#(K \setminus I)} \Pi_I.$$

The next result shows that Π_K^o describes an orthogonal projection, too. And the corresponding subspaces $\Pi_K^o \mathbb{H}$, $K \subset \{1, \dots, n\}$, comprise a decomposition of \mathbb{H} into pairwise orthogonal subspaces.

Theorem A.17. (a) For arbitrary sets $K \subset \{1, \dots, n\}$.

$$\Pi_K = \sum_{J \subset K} \Pi_J^o.$$

In particular, the identity operator I may be written as

$$I = \sum_{K \subset \{1, \dots, n\}} \Pi_K^o.$$

(b) For arbitrary sets $J, K \in \{1, \dots, n\}$,

$$\Pi_J^o \Pi_K = \Pi_K \Pi_J^o = 1_{[J \subset K]} \Pi_J^o,$$

and

$$\Pi_J^o \Pi_K^o = 1_{[J=K]} \Pi_J^o.$$

(c) Each operator Π_K^o describes the orthogonal projection of \mathbb{H} onto the linear space

$$\mathbb{H}_K^o := \mathbb{H}_K \cap \left(\sum_{J \subsetneq K} \mathbb{H}_J \right)^\perp.$$

These spaces \mathbb{H}_K^o , $K \subset \{1, \dots, n\}$, are pairwise orthogonal.

Corollary A.18 (Hoeffding's decomposition). *Any random variable $Y \in \mathbb{H}$ can be written as*

$$Y = \sum_{K \subset \{1, \dots, n\}} \Pi_K^\circ Y,$$

and the random variables $\Pi_K^\circ Y$, $K \subset \{1, \dots, n\}$, are uncorrelated with $\Pi_\emptyset^\circ Y \equiv \mathbb{E}(Y)$ and $\mathbb{E}(\Pi_K^\circ Y) = 0$ if $K \neq \emptyset$.

For an additional random variable $Z \in \mathbb{H}$,

$$\mathbb{E}(YZ) = \sum_{K \subset \{1, \dots, n\}} \mathbb{E}(\Pi_K^\circ Y \Pi_K^\circ Z).$$

This corollary follows essentially from Theorem A.17, except for the statements about $\Pi_\emptyset^\circ Y$ and $\mathbb{E}(\Pi_K^\circ Y)$. But $\mathbb{H}_\emptyset^\circ = \mathbb{H}_\emptyset$ is the space of constants, so $\Pi_\emptyset^\circ Y = \Pi_\emptyset Y \equiv \mathbb{E}(Y)$, and for any nonempty set $K \subset \{1, \dots, n\}$, it follows from $\mathbb{H}_K^\circ \perp \mathbb{H}_\emptyset$ that $\mathbb{E}(\Pi_K^\circ Y) = \langle \Pi_K^\circ Y, 1 \rangle = 0$.

Example A.19 (The case $n = 2$). Suppose that $Y = f(X_1, X_2)$. Then the Hoeffding decomposition of Y reads

$$Y - \mathbb{E}(Y) = f_1^\circ(X_1) + f_2^\circ(X_2) + f_{12}^\circ(X_1, X_2).$$

The three summands on the right hand side are given by

$$f_1^\circ(x_1) := \mathbb{E} f(x_1, X_2) - \mathbb{E}(Y), \quad f_2^\circ(x_2) := \mathbb{E} f(X_1, x_2) - \mathbb{E}(Y)$$

and

$$f_{12}^\circ(x_1, x_2) := f(x_1, x_2) - \mathbb{E} f(x_1, X_2) - \mathbb{E} f(X_1, x_2) + \mathbb{E}(Y).$$

Moreover, for arbitrary random variables $g_1(X_1)$ and $g_2(X_2)$ in $L^2(\mathbb{P})$,

$$\mathbb{E}(f_1^\circ(X_1)g_2(X_2)) = 0 = \mathbb{E}(f_{12}^\circ(X_1, X_2)g_2(X_2))$$

and

$$\mathbb{E}(f_2^\circ(X_2)g_1(X_1)) = 0 = \mathbb{E}(f_{12}^\circ(X_1, X_2)g_1(X_1)).$$

In particular, the random variables $f_1^\circ(X_1)$, $f_2^\circ(X_2)$ and $f_{12}^\circ(X_1, X_2)$ are centered and uncorrelated.

Proof of Theorem A.17. We start with a simple combinatorial fact. For any finite set S ,

$$(A.3) \quad \sum_{L \subset S} (-1)^{\#L} = 1_{[S=\emptyset]}.$$

This follows essentially from the binomial formula, for

$$\begin{aligned} \sum_{L \subset S} (-1)^{\#L} &= \sum_{\ell=0}^{\#S} \#\{L \subset S : \#L = \ell\} (-1)^\ell \\ &= \sum_{\ell=0}^{\#S} \binom{\#S}{\ell} (-1)^\ell (+1)^{\#S-\ell} = (1-1)^{\#S} = 1_{[S=\emptyset]}. \end{aligned}$$

As to part (a), by definition of Π_K^o and formula (A.3),

$$\begin{aligned} \sum_{J \subset K} \Pi_J^o &= \sum_{J \subset K} \sum_{I \subset J} (-1)^{\#(J \setminus I)} \Pi_I \\ &= \sum_{I \subset K} \left(\sum_{J \subset K: I \subset J} (-1)^{\#(J \setminus I)} \right) \Pi_I \\ &= \sum_{I \subset K} \left(\sum_{L \subset K \setminus I} (-1)^{\#L} \right) \Pi_I = \sum_{I \subset K} 1_{[K \setminus I = \emptyset]} \Pi_I = \Pi_K. \end{aligned}$$

As to part (b), it follows from (A.2) that

$$\begin{aligned} \left. \begin{array}{l} \Pi_J^o \Pi_K \\ \Pi_K \Pi_J^o \end{array} \right\} &= \sum_{I \subset J} (-1)^{\#(J \setminus I)} \Pi_{I \cap K} \\ &= \sum_{\tilde{I} \subset J \cap K} \sum_{L \subset J \setminus K} (-1)^{\#(J \setminus (\tilde{I} \cup L))} \Pi_{\tilde{I}} \\ &= \sum_{\tilde{I} \subset J \cap K} (-1)^{\#(J \setminus \tilde{I})} \left(\sum_{L \subset J \setminus K} (-1)^{-\#L} \right) \Pi_{\tilde{I}} \\ &= \sum_{\tilde{I} \subset J \cap K} (-1)^{\#(J \setminus \tilde{I})} 1_{[J \subset K]} \Pi_{\tilde{I}} = 1_{[J \subset K]} \sum_{\tilde{I} \subset J} (-1)^{\#(J \setminus \tilde{I})} \Pi_{\tilde{I}} = 1_{[J \subset K]} \Pi_J^o, \end{aligned}$$

where the second to last step follows from $(-1)^{-\#L} = (-1)^{\#L}$ and formula (A.3). This proves the first identities of part (b), and the second one follows from

$$\begin{aligned} \Pi_J^o \Pi_K^o &= \sum_{I \subset K} (-1)^{\#(K \setminus I)} \Pi_J^o \Pi_I \\ &= \sum_{I \subset K} (-1)^{\#(K \setminus I)} 1_{[J \subset I]} \Pi_J^o \\ &= 1_{[J \subset K]} \left(\sum_{I \subset K: J \subset I} (-1)^{\#(K \setminus I)} \right) \Pi_J^o \\ &= 1_{[J \subset K]} \left(\sum_{L \subset K \setminus J} (-1)^{\#L} \right) \Pi_J^o = 1_{[J \subset K]} 1_{[K \subset J]} \Pi_J^o = 1_{[J=K]} \Pi_J^o. \end{aligned}$$

It remains to prove part (c). As shown in Exercise A.20, a linear operator $\Pi : \mathbb{H} \rightarrow \mathbb{H}$ describes an orthogonal projection if and only if it satisfies $\Pi^2 = \Pi$ and is self-adjoint, that means $\langle \Pi Y, Z \rangle = \langle Y, \Pi Z \rangle$ for all $Y, Z \in \mathbb{H}$. By definition, all operators Π_J , $J \subset \{1, \dots, n\}$, have these properties, so Π_K^o , being a linear combination of self-adjoint operators, is self-adjoint, too. Moreover, it follows from part (b) that $\Pi_K^o \Pi_K^o = \Pi_K^o$, whence Π_K^o describes the orthogonal projection of \mathbb{H} onto some linear subspace \mathbb{H}_K^o . The subspaces \mathbb{H}_K^o , $K \subset \{1, \dots, n\}$, are pairwise orthogonal, because for different index sets J, K and $Y \in \mathbb{H}_J^o$, $Z \in \mathbb{H}_K^o$,

$$\langle Y, Z \rangle = \langle \Pi_J^o Y, \Pi_K^o Z \rangle = \langle Y, \Pi_J^o \Pi_K^o Z \rangle = \langle Y, 0 \rangle = 0.$$

Finally, by part (a), for any $K \subset \{1, \dots, n\}$,

$$\mathbb{H}_K = \sum_{J \subset K} \mathbb{H}_J^o = \mathbb{H}_K^o + \sum_{J \subsetneq K} \mathbb{H}_J^o,$$

so \mathbb{H}_K^o equals

$$\mathbb{H}_K \cap \left(\sum_{J \subsetneq K} \mathbb{H}_J^o \right)^\perp.$$

But $\mathbb{H}_I \subset \mathbb{H}_J$ for $I \subset J$, so $\mathbb{H}_J^o \subset \mathbb{H}_I$, whence

$$\sum_{J \subsetneq K} \mathbb{H}_J^o \subset \sum_{J \subsetneq K} \mathbb{H}_J.$$

On the other hand, for any fixed $\tilde{J} \subsetneq K$,

$$\sum_{J \subsetneq K} \mathbb{H}_J^o \supset \sum_{J \subset \tilde{J}} \mathbb{H}_J^o = \mathbb{H}_{\tilde{J}},$$

so

$$\sum_{J \subsetneq K} \mathbb{H}_J^o \supset \sum_{J \subsetneq K} \mathbb{H}_J.$$

Consequently,

$$\sum_{J \subsetneq K} \mathbb{H}_J^o = \sum_{J \subsetneq K} \mathbb{H}_J,$$

and this leads to the asserted representation of \mathbb{H}_K^o . \square

Exercise A.20 (Projections and orthogonal projections). Let $(\mathbb{H}, \langle \cdot, \cdot \rangle)$ be a real Hilbert space, and let $\Pi : \mathbb{H} \rightarrow \mathbb{H}$ be a linear mapping which is idempotent, that means, $\Pi^2 = \Pi$.

(a) Show that there exist linear subspaces $\mathbb{H}_1, \mathbb{H}_2$ of \mathbb{H} such that $\mathbb{H}_1 \cap \mathbb{H}_2 = \{0\}$, $\mathbb{H}_1 + \mathbb{H}_2 = \mathbb{H}$ and

$$\Pi x = \begin{cases} x & \text{if } x \in \mathbb{H}_1, \\ 0 & \text{if } x \in \mathbb{H}_2. \end{cases}$$

Hint: Write $x \in \mathbb{H}$ as $x = x_1 + x_2$ with $x_1 = \Pi x$ and $x_2 = x - \Pi x$.

(b) Show that $\mathbb{H}_1 \perp \mathbb{H}_2$ if and only if Π is self-adjoint, that means, $\langle \Pi x, y \rangle = \langle x, \Pi y \rangle$ for all $x, y \in \mathbb{H}$. In this case, Π is the orthogonal projection onto \mathbb{H}_1 .

Exercise A.21. Let $(\mathbb{H}, \langle \cdot, \cdot \rangle)$ be a real Hilbert space, and let Π_1, Π_2 be orthogonal projections onto subspaces \mathbb{H}_1 and \mathbb{H}_2 , respectively. Further let Π_0 be the orthogonal projection onto $\mathbb{H}_0 := \mathbb{H}_1 \cap \mathbb{H}_2$. Show that the following three statements are equivalent:

(i) $\mathbb{H}_1 \cap \mathbb{H}_0^\perp \perp \mathbb{H}_2 \cap \mathbb{H}_0^\perp$.

(ii) $\Pi_1 \Pi_2 = \Pi_0$.

(iii) $\Pi_1 \Pi_2 = \Pi_2 \Pi_1$.

A.5 Weak Law of Large Numbers and Central Limit Theorem

In connection with asymptotic considerations, the subsequent versions of the Weak Law of Large Numbers and Lindeberg's Central Limit Theorem are rather useful. Throughout this section asymptotic statements refer to $n \rightarrow \infty$, unless specified differently.

Theorem A.22 (WLLN). For any integer $n \geq 1$ let $Y_{n1}, Y_{n2}, \dots, Y_{nn}$ be independent random variables such that

$$\begin{aligned} \sum_{i=1}^n \mathbb{E} |Y_{ni}| &= O(1), \\ \sum_{i=1}^n \mathbb{E} (1_{[|Y_{ni}| > \epsilon]} |Y_{ni}|) &\rightarrow 0 \quad \text{for any fixed } \epsilon > 0. \end{aligned}$$

Then with $\mu_{ni} := \mathbb{E}(Y_{ni})$,

$$\mathbb{E} \left| \sum_{i=1}^n (Y_{ni} - \mu_{ni}) \right| \rightarrow 0 \quad \text{and} \quad \mathbb{E} \max_{1 \leq i \leq n} |Y_{ni}| \rightarrow 0.$$

Theorem A.23 (CLT). For any integer $n \geq 1$ let $Y_{n1}, Y_{n2}, \dots, Y_{nn}$ be independent random variables such that for some real numbers μ and $\sigma > 0$,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(Y_{ni}) &\rightarrow \mu \quad \text{and} \quad \sum_{i=1}^n |\mathbb{E}(Y_{ni})| = O(1), \\ \sum_{i=1}^n \mathbb{E}(Y_{ni}^2) &\rightarrow \sigma^2, \\ \sum_{i=1}^n \mathbb{E}(1_{[Y_{ni}^2 > \epsilon]} Y_{ni}^2) &\rightarrow 0 \quad \text{for any fixed } \epsilon > 0. \end{aligned}$$

Then

$$\sum_{i=1}^n Y_{ni} \rightarrow_{\mathcal{L}} \mathcal{N}(\mu, \sigma^2)$$

and

$$\mathbb{E} \left| \sum_{i=1}^n Y_{ni}^2 - \sigma^2 \right| \rightarrow 0, \quad \mathbb{E} \left(\max_{1 \leq i \leq n} Y_{ni}^2 \right) \rightarrow 0, \quad \sum_{i=1}^n \mathbb{E}(Y_{ni})^2 \rightarrow 0.$$

Corollary A.24. For any integer $n \geq 1$ let $X_{n1}, X_{n2}, \dots, X_{nn}$ be independent and identically distributed random variables such that for some real numbers μ and $\sigma > 0$,

$$\begin{aligned} \sqrt{n} \mathbb{E}(X_{n1}) &\rightarrow \mu, \\ \mathbb{E}(X_{n1}^2) &\rightarrow \sigma^2, \\ \mathbb{E}(1_{[X_{n1}^2 > \epsilon n]} X_{n1}^2) &\rightarrow 0 \quad \text{for any fixed } \epsilon > 0. \end{aligned}$$

Then

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{ni} \rightarrow_{\mathcal{L}} \mathcal{N}(\mu, \sigma^2)$$

and

$$\mathbb{E} \left| \frac{1}{n} \sum_{i=1}^n X_{ni}^2 - \sigma^2 \right| \rightarrow 0, \quad \mathbb{E} \left(\max_{1 \leq i \leq n} \frac{X_{ni}^2}{n} \right) \rightarrow 0.$$

Proof of Theorem A.22. Let $M := \limsup_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{E} |Y_{ni}|$. For arbitrary fixed $\epsilon > 0$ set

$$Y_{ni1} := 1_{[|Y_{ni}| \leq \epsilon]} Y_{ni}, \quad Y_{ni2} := 1_{[|Y_{ni}| > \epsilon]} Y_{ni}$$

and $\mu_{nik} := \mathbb{E}(Y_{nik})$. Then

$$\begin{aligned}
\mathbb{E} \left| \sum_{i=1}^n (Y_{ni} - \mu_{ni}) \right| &\leq \mathbb{E} \left| \sum_{i=1}^n (Y_{ni1} - \mu_{ni1}) \right| + \sum_{i=1}^n (\mathbb{E} |Y_{ni2}| + |\mu_{ni2}|) \\
&\leq \sqrt{\text{Var} \left(\sum_{i=1}^n Y_{ni1} \right)} + 2 \sum_{i=1}^n \mathbb{E} |Y_{ni2}| \\
&\leq \sqrt{\sum_{i=1}^n \mathbb{E}(Y_{ni1}^2)} + o(1) \\
&\leq \sqrt{\epsilon \sum_{i=1}^n \mathbb{E} |Y_{ni}|} + o(1) \\
&\leq \sqrt{\epsilon M} + o(1) + o(1) \\
&\rightarrow \sqrt{\epsilon M}.
\end{aligned}$$

Furthermore,

$$\mathbb{E} \left(\max_{1 \leq i \leq n} |Y_{ni}| \right) \leq \epsilon + \sum_{i=1}^n \mathbb{E} |Y_{ni2}| \rightarrow \epsilon.$$

Since $\epsilon > 0$ may be arbitrarily small, these calculations yield the assertions. \square

Remarks on the proof of Theorem A.23. One can deduce from Theorem A.22, applied to Y_{ni}^2 in place of Y_{ni} , that

$$\mathbb{E} \left| \sum_{i=1}^n Y_{ni}^2 - \sigma^2 \right| \rightarrow 0 \quad \text{and} \quad \mathbb{E} \left(\max_{1 \leq i \leq n} Y_{ni}^2 \right) \rightarrow 0.$$

In particular, with $\mu_{ni} := \mathbb{E}(Y_{ni})$,

$$\max_{1 \leq i \leq n} \mu_{ni}^2 \leq \max_{1 \leq i \leq n} \mathbb{E}(Y_{ni}^2) \leq \mathbb{E} \left(\max_{1 \leq i \leq n} Y_{ni}^2 \right) \rightarrow 0,$$

whence

$$\sum_{i=1}^n \mu_{ni}^2 \leq \max_{1 \leq j \leq n} |\mu_{nj}| \sum_{i=1}^n |\mu_{ni}| \rightarrow 0.$$

But with $\sigma_{ni}^2 := \text{Var}(Y_{ni})$ this implies that

$$\sum_{i=1}^n \sigma_{ni}^2 = \sum_{i=1}^n \mathbb{E}(Y_{ni}^2) - \sum_{i=1}^n \mu_{ni}^2 \rightarrow \sigma^2.$$

Moreover, for any fixed $\epsilon > 0$, the inequality $\max_{1 \leq i \leq n} |\mu_{ni}| \leq \epsilon/2$ is satisfied for sufficiently large n , and in that case, $|Y_{ni} - \mu_{ni}| > \epsilon$ implies that $|Y_{ni}| \geq \epsilon/2$ and $|Y_{ni} - \mu_{ni}| \leq 2|Y_{ni}|$. Hence for sufficiently large n ,

$$\sum_{i=1}^n \mathbb{E} (1_{\{|Y_{ni} - \mu_{ni}| > \epsilon\}} (Y_{ni} - \mu_{ni})^2) \leq 4 \sum_{i=1}^n \mathbb{E} (1_{\{|Y_{ni}| > \epsilon/2\}} Y_{ni}^2) \rightarrow 0.$$

Consequently, the centered random variables $Z_{ni} := Y_{ni} - \mu_{ni}$ satisfy the assumptions of the more traditional CLT:

$$\begin{aligned} \mathbb{E}(Z_{ni}) &= 0 \quad \text{for all } n \geq 1 \text{ and } 1 \leq i \leq n, \\ \sum_{i=1}^n \mathbb{E}(Z_{ni}^2) &\rightarrow \sigma^2, \\ \sum_{i=1}^n \mathbb{E}(1_{\{|Z_{ni}| > \epsilon\}} Z_{ni}^2) &\rightarrow 0 \quad \text{for any fixed } \epsilon > 0. \end{aligned}$$

These conditions imply that

$$\sum_{i=1}^n Z_{ni} \rightarrow_{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

and thus

$$\sum_{i=1}^n Y_{ni} = \mu + o(1) + \sum_{i=1}^n Z_{ni} \rightarrow_{\mathcal{L}} \mathcal{N}(\mu, \sigma^2).$$

□