

AN INTRODUCTION TO ESTIMATION AND DETECTION THEORY¹

Armand M. Makowski² and Prakash L.K. Narayan³

¹©1994-2017 by Armand M. Makowski and Prakash L.K. Narayan

²Department of Electrical and Computer Engineering, and Institute for Systems Research, University of Maryland, College Park, MD 20742. E-mail: armand@isr.umd.edu. Phone: (301) 405-6844

³Department of Electrical and Computer Engineering, and Institute for Systems Research, University of Maryland, College Park, MD 20742. E-mail: prakash@isr.umd.edu. Phone: (301) 405-3661

Notation and conventions

Throughout, we use \mathbb{R} to denote the set of all real numbers, or equivalently, the real line $(-\infty, \infty)$. The set of all non-negative real numbers is denoted by \mathbb{R}_+ . The set of $\{0, 1, \dots\}$ of all non-negative integers is denoted \mathbb{N} . The notation \mathbb{N}_0 will be used for the set $\{1, 2, \dots\}$ of all positive integers.

With p a positive integer, let \mathbb{R}^p denote the p^{th} cartesian product of \mathbb{R} . An element \boldsymbol{x} in \mathbb{R}^p , whose p components are denoted x_1, \dots, x_p , is always interpreted as a column vector $(x_1, \dots, x_p)'$ (with $'$ denoting transpose).

Part I
DETECTION THEORY

Chapter 1

Simple binary hypothesis testing

A decision has to be made as to which of two hypotheses (or states of nature) is the correct one. The states of nature are encoded in a rv H and a decision has to be made on the basis of an observation Y which is statistically related to H .

1.1 Motivating examples

Control process A machine produces circuit boards. It is either fully functioning ($H = 1$) or worn out ($H = 0$). Checking the state of the machine is not feasible as it would require that the production be stopped, incurring a loss of revenue for the manufacturer if the machine were indeed shown to be fully functionally. Instead, a batch of circuits is collected and tested for a number of performance parameters, say Y_1, \dots, Y_k . It is known that

A simple communication example

Testing means

1.2 The probabilistic model

These examples can be cast as *binary* hypothesis testing problems: Nature is in either of two states, say $H = 0$ or $H = 1$ for sake of concreteness, and the observations are organized into an \mathbb{R}^k -valued rv \mathbf{Y} . We assume given two probability distribution functions $F_0, F_1 : \mathbb{R}^k \rightarrow [0, 1]$ on \mathbb{R}^k ; they will act as conditional

probability distribution of \mathbf{Y} given $H = 0$ and $H = 1$, respectively. This situation is summarized by

$$\begin{aligned} H_1 : Y &\sim F_1 \\ H_0 : Y &\sim F_0. \end{aligned} \quad (1.1)$$

In the statistical literature the hypothesis H_0 is called the *null hypothesis* and hypothesis H_1 is referred to as the *non-null hypothesis* or the *alternative*.

Probabilistically, the symbolic statement (1.1) is understood as follows: Given some probability triple $(\Omega, \mathcal{F}, \mathbb{P})$ (whose existence is discussed shortly), consider rvs $H : \Omega \rightarrow \{0, 1\}$ and $\mathbf{Y} : \Omega \rightarrow \mathbb{R}^k$. The probability distribution functions F_0 and F_1 being interpreted as conditional probability distribution of \mathbf{Y} given $H = 0$ and $H = 1$, respectively, we must have

$$F_h(\mathbf{y}) = \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = h], \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ h = 0, 1. \end{array}$$

The probability distribution of the rv H is specified by p in $[0, 1]$ with

$$p = \mathbb{P}[H = 1] = 1 - \mathbb{P}[H = 0].$$

We refer to the pmf $(1 - p, p)$ on $\{0, 1\}$, or just to p , as the *prior*.

Because

$$\begin{aligned} \mathbb{P}[\mathbf{Y} \leq \mathbf{y}, H = h] &= \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = h] \mathbb{P}[H = h] \\ &= \begin{cases} (1 - p)F_0(\mathbf{y}) & \text{if } h = 0, \mathbf{y} \in \mathbb{R}^k \\ pF_1(\mathbf{y}) & \text{if } h = 1, \mathbf{y} \in \mathbb{R}^k, \end{cases} \end{aligned} \quad (1.2)$$

the law of total probability shows that

$$\begin{aligned} \mathbb{P}[\mathbf{Y} \leq \mathbf{y}] &= \sum_{h=0}^1 \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = h] \mathbb{P}[H = h] \\ &= pF_1(\mathbf{y}) + (1 - p)F_0(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned} \quad (1.3)$$

In other words, the conditional probability distributions of the observations given the hypothesis *and* the probability distribution of H completely specify the *joint* distribution of the rvs H and \mathbf{Y} .

1.3 A construction

The existence of the model described in Section 1.2 can be guaranteed through the following construction: Take $\Omega = \{0, 1\} \times \mathbb{R}^k$ with generic element $\omega = (h, \mathbf{y})$ with $h = 0, 1$ and \mathbf{y} an arbitrary element of \mathbb{R}^k . We endow Ω with the σ -field \mathcal{F} given by

$$\mathcal{F} = \sigma(\mathcal{P}(\{0, 1\}) \times \mathcal{B}(\mathbb{R}^k))$$

where $\mathcal{P}(\{0, 1\})$ is the power set of $\{0, 1\}$, and $\mathcal{B}(\mathbb{R}^k)$ is the Borel σ -field on \mathbb{R}^k .

We define the mappings $H : \Omega \rightarrow \mathbb{R}$ and $\mathbf{Y} : \Omega \rightarrow \mathbb{R}^k$ by

$$H(\omega) = h \quad \text{and} \quad \mathbf{Y}(\omega) = \mathbf{y}, \quad \omega = (h, \mathbf{y}) \in \Omega.$$

Both projection mappings are Borel measurable, and therefore define rvs.

If \mathbb{P} is any probability measure on the σ -field \mathcal{F} , then by *construction* of the rvs H and \mathbf{Y} just given, the joint probability distribution of the pair (H, \mathbf{Y}) is necessarily given by

$$\begin{aligned} \mathbb{P}[H = h, \mathbf{Y} \leq \mathbf{y}] &= \mathbb{P}[\{\omega \in \Omega : H(\omega) = h, \mathbf{Y}(\omega) \leq \mathbf{y}\}] \\ &= \mathbb{P}[\{h\} \times (-\infty, \mathbf{y}]], \quad \begin{array}{l} h = 0, 1 \\ \mathbf{y} \in \mathbb{R}^k \end{array} \end{aligned} \quad (1.4)$$

since

$$\{\omega \in \Omega : H(\omega) = h, \mathbf{Y}(\omega) \leq \mathbf{y}\} = \{h\} \times (-\infty, \mathbf{y}].$$

On the way to identify a probability \mathbb{P} on \mathcal{F} under which the joint probability distribution of the pair (H, \mathbf{Y}) satisfies (1.2), we readily conclude from (1.4) that \mathbb{P} is necessarily determined on certain rectangles, namely

$$\mathbb{P}[\{h\} \times (-\infty, \mathbf{y}]] = \begin{cases} (1-p)F_0(\mathbf{y}) & \text{if } h = 0 \\ pF_1(\mathbf{y}) & \text{if } h = 1 \end{cases} \quad (1.5)$$

for every \mathbf{y} in \mathbb{R}^k . At this point we recall the following fact from Measure Theory: Any probability measure on the σ -field \mathcal{F} carried by the product space $\{0, 1\} \times \mathbb{R}^k$ is uniquely determined on the entire σ -field \mathcal{F} by its values on the rectangle sets of the form

$$\{h\} \times (-\infty, \mathbf{y}], \quad \begin{array}{l} h = 0, 1 \\ \mathbf{y} \in \mathbb{R}^k. \end{array}$$

Therefore, by virtue of (1.5) there exists a *unique* probability measure \mathbb{P} on \mathcal{F} such that (1.4) holds. More generally, it is also the case that

$$\mathbb{P}[\{h\} \times (-\infty, \mathbf{y}]] = \begin{cases} (1-p) \int_B dF_0(\mathbf{y}) & \text{if } h = 0 \\ p \int_B dF_1(\mathbf{y}) & \text{if } h = 1 \end{cases} \quad (1.6)$$

for every Borel set B in \mathbb{R}^k as a result of the fact that

$$\sigma((-\infty, \mathbf{y}], \mathbf{y} \in \mathbb{R}^k) = \mathcal{B}(\mathbb{R}^k).$$

Finally, under this probability measure \mathbb{P} it is plain (1.5) immediately implies

$$\mathbb{P}[H = h] = \mathbb{P}[\{h\}] = \begin{cases} (1-p) & \text{if } h = 0 \\ p & \text{if } h = 1 \end{cases} \quad (1.7)$$

and

$$\begin{aligned} \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = h] &= \frac{\mathbb{P}[H = h, \mathbf{Y} \leq \mathbf{y}]}{\mathbb{P}[H = h]} \\ &= F_h(\mathbf{y}) \end{aligned} \quad (1.8)$$

for every \mathbf{y} in \mathbb{R}^k , as required.

1.4 Basic assumptions

During the discussion, several assumptions will be enforced on the probability distributions F_0 and F_1 . The assumptions that will be most often encountered are denoted by **(A.1)** and **(A.2)** for sake of convenience. They are stated and discussed in some details below.

Condition (A.1): The probability distributions F_0 and F_1 on \mathbb{R}^k are both *absolutely continuous* with respect to some distribution F on \mathbb{R}^k – In general F may not be a probability distribution.

Condition (A.1) is equivalent to saying that there exist Borel mappings $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ such that

$$F_h(\mathbf{y}) = \int_{-\infty}^{\mathbf{y}} f_h(\boldsymbol{\eta}) dF(\boldsymbol{\eta}), \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ h = 0, 1. \end{array} \quad (1.9)$$

In some basic sense, this condition is hardly constraining since we can always take F to be the average of the two probability distributions F_0 and F_1 . i.e.,

$$F(\mathbf{y}) \equiv \frac{1}{2}F_0(\mathbf{y}) + \frac{1}{2}F_1(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \quad (1.10)$$

in which case F is also a probability distribution. This choice for F is usually not operationally convenient and therefore discarded. However, the most often encountered situations arise when F is either Lebesgue measure on \mathbb{R}^k or a counting measure on some countable subset of \mathbb{R}^k , in which case F is not a probability distribution.

When F is Lebesgue measure on \mathbb{R}^k , the Borel mappings $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ are just the probability density functions induced by F_0 and F_1 in the usual sense. When F is counting measure on a countable subset $S \subseteq \mathbb{R}^k$, then the Borel mappings $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ are best thought as *probability mass functions* (pdfs) $\mathbf{f}_0 = \{f_0(\mathbf{y}), \mathbf{y} \in S\}$ and $\mathbf{f}_1 = \{f_1(\mathbf{y}), \mathbf{y} \in S\}$, i.e.,

$$0 \leq f_h(\mathbf{y}) \leq 1, \quad \begin{array}{l} \mathbf{y} \in S, \\ h = 0, 1. \end{array}$$

and

$$\sum_{\mathbf{y} \in S} f_h(\mathbf{y}) = 1, \quad h = 0, 1.$$

The condition (1.9) now takes the form

$$\mathbb{P}[\mathbf{Y} \in B | H = h] = \sum_{\boldsymbol{\eta} \in S \cap B} f_h(\boldsymbol{\eta}), \quad \begin{array}{l} B \in \mathcal{B}(\mathbb{R}^k) \\ h = 0, 1. \end{array}$$

Condition **(A.2)**: The probability distribution F_1 is *absolutely continuous* with respect to the probability distribution F_0 .

Under Condition **(A.1)**, with the notation introduced earlier, this is equivalent to requiring

$$f_0(\mathbf{y}) = 0 \quad \text{implies} \quad f_1(\mathbf{y}) = 0. \quad (1.11)$$

1.5 Admissible tests

Decisions as to which state of nature occurred are taken on the basis of observations; this is formalized through the following definition.

An *admissible* decision rule (or test) is any *Borel* mapping $d : \mathbb{R}^k \rightarrow \{0, 1\}$. The collection of all admissible rules is denoted by \mathcal{D} .

The measurability requirement entering the definition of admissibility is imposed to guarantee that the mapping $d(\mathbf{Y}) : \Omega \rightarrow \{0, 1\} : \omega \rightarrow d(Y(\omega))$ is indeed a rv, i. e., $[\omega \in \Omega : d(Y(\omega)) = h]$ is an event in \mathcal{F} for all $h = 0, 1$. The need for this technical condition will become apparent in subsequent chapters.

The next fact will prove useful in some of the discussion

Lemma 1.5.1 *The set \mathcal{D} of admissible decision rules is in one-to-one correspondence with $\mathcal{B}(\mathbb{R}^k)$.*

Proof. By definition of admissibility every test d in \mathcal{D} is completely specified by the *Borel* subset $C(d)$ defined by

$$C(d) \equiv \{\mathbf{y} \in \mathbb{R}^k : d(\mathbf{y}) = 0\}. \quad (1.12)$$

Conversely, any Borel measurable subset C of \mathbb{R}^k uniquely determines an admissible rule d_C in \mathcal{D} through

$$d_C(\mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{y} \notin C \\ 0 & \text{if } \mathbf{y} \in C. \end{cases}$$

We note that $C(d_C) = C$ as expected. ■

Any admissible rule d in \mathcal{D} induces *two* types of error: Upon observing \mathbf{Y} , either $H = 0$ is true and $d(\mathbf{Y}) = 1$ or $H = 1$ is true and $d(\mathbf{Y}) = 0$.

These two possibilities are the so-called errors of the *first* and *second* type associated with the decision rule d ; they are quantified by

$$\alpha(d) \equiv \mathbb{P}[d(\mathbf{Y}) = 1 | H = 0] \quad (1.13)$$

and

$$\beta(d) \equiv \mathbb{P}[d(\mathbf{Y}) = 0 | H = 1], \quad (1.14)$$

respectively.

The quantity $\alpha(d)$ is sometimes called the *size* of the test d . In radar parlance, these probabilities are referred to as probabilities of *false alarm* and *miss*, respectively, with alternate notation

$$P_F(d) \equiv \mathbb{P}[d(\mathbf{Y}) = 1 | H = 0] \quad (1.15)$$

and

$$P_M(d) \equiv \mathbb{P}[d(\mathbf{Y}) = 0 | H = 1]. \quad (1.16)$$

Throughout we shall use this terminology. Sometimes, it is convenient to consider the so-called probability of *detection* given by

$$P_D(d) \equiv \mathbb{P}[d(\mathbf{Y}) = 1 | H = 1] = 1 - P_M(d). \quad (1.17)$$

1.6 Likelihood ratio tests

In subsequent chapters we shall consider several formulations for the binary hypothesis problem. In all cases the tests of interest are related to tests in the class of admissible tests $\{d_\eta, \eta \geq 0\}$ which we now introduce.

For each $\eta \geq 0$, the mapping $d_\eta : \mathbb{R}^k \rightarrow \{0, 1\}$ is defined by

$$d_\eta(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < \eta f_0(\mathbf{y}). \quad (1.18)$$

It is plain from the definition (1.18) (with $\eta = 0$) that d_0 is simply the test that always selects the non-null hypothesis $H = 1$, i.e., $d_0(\mathbf{y}) = 1$ for every \mathbf{y} in \mathbb{R}^k . On the other hand, formally substituting $\eta = \infty$ in (1.18) will be problematic at observation points where $f_0(\mathbf{y}) = 0$. However, by *convention* we shall interpret d_∞ as the test that always selects the null hypothesis $H = 0$, i.e., $d_\infty(\mathbf{y}) = 0$ for every \mathbf{y} in \mathbb{R}^k .

Such tests take an even simpler form under the additional Condition **(A.2)** as will be seen shortly: Note that (1.18) can be rewritten as

$$d_\eta(\mathbf{y}) = 0 \quad \text{if} \quad \frac{f_1(\mathbf{y})}{f_0(\mathbf{y})} < \eta \quad \text{whenever} \quad f_0(\mathbf{y}) > 0.$$

Taking our cue from this last statement, we define the *likelihood ratio* as any Borel mapping $L : \mathbb{R}^k \rightarrow \mathbb{R}$ of the form

$$L(\mathbf{y}) \equiv \begin{cases} \frac{f_1(\mathbf{y})}{f_0(\mathbf{y})} & \text{if } f_0(\mathbf{y}) > 0 \\ \Lambda(\mathbf{y}) & \text{if } f_0(\mathbf{y}) = 0 \end{cases} \quad (1.19)$$

for some arbitrary Borel mapping $\Lambda : \mathbb{R}^k \rightarrow \mathbb{R}_+$. Different choices of this arbitrary non-negative function produce different versions of the likelihood ratio function.

Given a version of the likelihood ratio function in (1.19), we define the *likelihood ratio* test with *threshold* $\eta \geq 0$ to be the admissible decision rule $Lrt_\eta : \mathbb{R}^k \rightarrow \{0, 1\}$ given by

$$Lrt_\eta(\mathbf{y}) \equiv \begin{cases} 1 & \text{if } L(\mathbf{y}) \geq \eta \\ 0 & \text{if } L(\mathbf{y}) < \eta. \end{cases} \quad (1.20)$$

With

$$B_h = \{\mathbf{y} \in \mathbb{R}^k : f_h(\mathbf{y}) = 0\}, \quad h = 0, 1, \quad (1.21)$$

we note that

$$\mathbb{P}[f_0(\mathbf{Y}) = 0 | H = h] = \int_{B_0} f_h(\mathbf{y}) dF(\mathbf{y}), \quad h = 0, 1. \quad (1.22)$$

Under **(A.2)**, the inclusion $B_0 \subseteq B_1$ holds and we conclude that

$$\mathbb{P}[f_0(\mathbf{Y}) = 0 | H = h] = 0, \quad h = 0, 1.$$

For any value η of the threshold it is plain that the tests d_η and Lrt_η coincide on the set $\{\mathbf{y} \in \mathbb{R}^k : f_0(\mathbf{y}) > 0\}$ (while possibly disagreeing on the complement B_0). Thus, for each $h = 0, 1$, we find that

$$\begin{aligned} & \mathbb{P}[d_\eta(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P}[d_\eta(\mathbf{Y}) = 0, f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P}[d_\eta(\mathbf{Y}) = 0, f_0(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 0, f_0(\mathbf{Y}) > 0 | H = h] \\ &= \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 0, f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 0, f_0(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 0 | H = h]. \end{aligned}$$

This discussion leads to the following fact.

Lemma 1.6.1 *Assume the absolute continuity conditions (A.1)–(A.2) to hold. For each $\eta \geq 0$, the tests d_η and Lrt_η are equivalent in the sense that $P_M(d_\eta) = P_M(Lrt_\eta)$ and $P_F(d_\eta) = P_F(Lrt_\eta)$.*

1.7 Exercises

1.8 References

Chapter 2

The Bayesian formulation

The Bayesian formulation assumes *knowledge* of the conditional distributions F_1 and F_0 , and of the prior distribution p of the rv H . Two other formulations, namely the Minimax formulation and the Neyman-Pearson formulation, will be studied in Chapters 4 and 5, respectively.

2.1 The Bayesian optimization problem

The cost incurred for making decisions is quantified by the mapping $C : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$ with the interpretation that

$$C(h, d) = \begin{array}{l} \text{Cost incurred for deciding } d \\ \text{when } H = h \end{array}, \quad d, h = 0, 1.$$

As the sample ω in Ω is realized, the observation $Y(\omega)$ is recorded and the use of the admissible rule d in \mathcal{D} incurs a cost $C(H(\omega), d(\mathbf{Y}(\omega)))$. Although it is tempting to seek to minimize this quantity, this is not possible. Indeed, the rv \mathbf{Y} is observed, whence $d(\mathbf{Y})$ is known once the test d has been specified, but the state of nature H is *not* directly observable. Consequently, the value of the cost $C(H, d(\mathbf{Y}))$ is not available. To remedy to this difficulty, we introduce the *expected cost function* $J : \mathcal{D} \rightarrow \mathbb{R}$ given by

$$J(d) \equiv \mathbb{E}[C(H, d(\mathbf{Y}))], \quad d \in \mathcal{D}.$$

The *Bayesian Problem* \mathcal{P}_B is the minimization problem

$$\mathcal{P}_B : \quad \text{Minimize } J(d) \text{ over } d \text{ in } \mathcal{D}.$$

This amounts to finding an admissible test $d^* : \mathbb{R}^k \rightarrow \{0, 1, \dots, M - 1\}$ in \mathcal{D} such that

$$J(d^*) \leq J(d), \quad d \in \mathcal{D}. \quad (2.1)$$

Any admissible test d^* which satisfies (2.1) is called a Bayesian test, and the value

$$J(d^*) = \inf_{d \in \mathcal{D}} J(d) = \min_{d \in \mathcal{D}} J(d) \quad (2.2)$$

is known as the *Bayesian cost*.

The solution to the Bayesian problem \mathcal{P}_B is developed with the help of an auxiliary result concerning the form of the Bayesian cost. This representation result will be useful in several places and is given here for sake of easy reference: Introduce the *relative costs* Γ_0 and Γ_1 given by

$$\Gamma_h \equiv C(h, 1 - h) - C(h, h), \quad h = 0, 1 \quad (2.3)$$

and define the auxiliary expected cost function $\hat{J} : \mathcal{D} \rightarrow \mathbb{R}$ to be

$$\hat{J}(d) = \mathbb{E} [\mathbf{1} [d(\mathbf{Y}) \neq H] \Gamma_H], \quad d \in \mathcal{D}. \quad (2.4)$$

Lemma 2.1.1 *For any admissible rule d in \mathcal{D} , the relation*

$$J(d) = \mathbb{E} [C(H, H)] + \hat{J}(d) \quad (2.5)$$

holds with

$$\hat{J}(d) = \Gamma_0(1 - p) \cdot P_F(d) + \Gamma_1 p \cdot P_M(d). \quad (2.6)$$

Proof. Fix d in \mathcal{D} . Recall that the rvs H and $d(\mathbf{Y})$ are $\{0, 1\}$ -valued rvs, and that the events $[d(\mathbf{Y}) = H]$ and $[d(\mathbf{Y}) \neq H]$ form a partition of Ω , i.e.,

$$\mathbf{1} [d(\mathbf{Y}) = H] + \mathbf{1} [d(\mathbf{Y}) \neq H] = \mathbf{1} [\Omega] = 1.$$

It readily follows that

$$\begin{aligned}
C(H, d(\mathbf{Y})) &= \mathbf{1}[d(\mathbf{Y}) = H] C(H, H) + \mathbf{1}[d(\mathbf{Y}) \neq H] C(H, 1 - H) \\
&= (1 - \mathbf{1}[d(\mathbf{Y}) \neq H]) C(H, H) + \mathbf{1}[d(\mathbf{Y}) \neq H] C(H, 1 - H) \\
&= C(H, H) + (C(H, 1 - H) - C(H, H)) \mathbf{1}[d(\mathbf{Y}) \neq H] \\
&= C(H, H) + \mathbf{1}[d(\mathbf{Y}) \neq H] \Gamma_H
\end{aligned} \tag{2.7}$$

with the relative costs Γ_0 and Γ_1 given by (2.3). Taking expectations on both sides of (2.7) we obtain (2.5).

The law of total probabilities gives

$$\begin{aligned}
\hat{J}(d) &= \mathbb{E}[\Gamma_0 \mathbf{1}[d(\mathbf{Y}) \neq 0] \mathbf{1}[H = 0] + \Gamma_1 \mathbf{1}[d(\mathbf{Y}) \neq 1] \mathbf{1}[H = 1]] \\
&= \Gamma_0(1 - p) \cdot \mathbb{P}[d(\mathbf{Y}) \neq 0|H = 0] + \Gamma_1 p \cdot \mathbb{P}[d(\mathbf{Y}) \neq 1|H = 1] \\
&= \Gamma_0(1 - p) \cdot \mathbb{P}[d(\mathbf{Y}) = 1|H = 0] + \Gamma_1 p \cdot \mathbb{P}[d(\mathbf{Y}) = 0|H = 1],
\end{aligned}$$

and the desired expression (2.6) is obtained. \blacksquare

The Bayesian cost under a given decision rule is completely determined by its probabilities of false alarm and of miss. We also note that

$$\begin{aligned}
\hat{J}(d) &= \Gamma_0(1 - p) + \Gamma_1 p \cdot \mathbb{P}[d(\mathbf{Y}) = 0|H = 1] \\
&\quad - \Gamma_0(1 - p) \cdot \mathbb{P}[d(\mathbf{Y}) = 0|H = 0], \quad d \in \mathcal{D}
\end{aligned} \tag{2.8}$$

as an immediate consequence of (2.6).

Therefore, by Lemma 1.6.1 it follows from (2.5)-(2.6) that $J(d_\eta) = J(Lrt_\eta)$ regardless of the cost function $C : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$. The same argument also shows that any two versions of the likelihood ratio function will generate likelihood ratio tests which are equivalent.

2.2 Solving the Bayesian problem \mathcal{P}_B

It follows from (2.6) that solving \mathcal{P}_B is *equivalent* to solving the auxiliary problem $\hat{\mathcal{P}}_B$ where

$$\hat{\mathcal{P}}_B : \quad \text{Minimize } \hat{J}(d) \text{ over } d \text{ in } \mathcal{D}.$$

To do so, it will be necessary to assume that the probability distributions F_0 and F_1 satisfy the absolute continuity condition **(A1)** given earlier, namely that

there exists a single distribution F on \mathbb{R}^k with respect to which both F_0 and F_1 are absolutely continuous. For any test d in \mathcal{D} , we get

$$\begin{aligned}\mathbb{P}[d(\mathbf{Y}) = 0 | H = h] &= \int_{C(d)} dF_h(\mathbf{y}) \\ &= \int_{C(d)} f_h(\mathbf{y}) dF(\mathbf{y}), \quad h = 0, 1\end{aligned}\quad (2.9)$$

with $C(d)$ defined at (1.12). It is now easy to see from (2.8) that

$$\widehat{J}(d) = \Gamma_0(1 - p) + \int_{C(d)} h(\mathbf{y}) dF(\mathbf{y}) \quad (2.10)$$

where the mapping $h : \mathbb{R}^k \rightarrow \mathbb{R}$ is given by

$$h(\mathbf{y}) \equiv \Gamma_1 p \cdot f_1(\mathbf{y}) - \Gamma_0(1 - p) \cdot f_0(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \quad (2.11)$$

Theorem 2.2.1 *Assume the absolute continuity condition (A.1) to hold. Define the Borel set C^* by*

$$C^* \equiv \{y \in \mathbb{R}^k : h(\mathbf{y}) < 0\} \quad (2.12)$$

with $h : \mathbb{R}^k \rightarrow \mathbb{R}$ given by (2.11). The decision rule $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ induced by C^* is given by

$$d^*(\mathbf{y}) = \begin{cases} 1 & \text{if } x \notin C^* \\ 0 & \text{if } x \in C^*; \end{cases} \quad (2.13)$$

it is admissible and solves the Problem $\widehat{\mathcal{P}}_B$, hence solves the Bayesian Problem \mathcal{P}_B .

Proof. The set C^* is a Borel subset of \mathbb{R}^k due to the fact that the functions $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ are themselves Borel measurable. The test d^* is therefore an admissible decision rule in \mathcal{D} since $C(d^*) = C^*$. We now show that d^* satisfies

$$\widehat{J}(d^*) \leq \widehat{J}(d), \quad d \in \mathcal{D}. \quad (2.14)$$

Indeed, for every test d in \mathcal{D} , we see from (2.10) that

$$\widehat{J}(d) = \Gamma_0(1 - p) + \int_{C(d) \setminus C^*} h(\mathbf{y}) dF(\mathbf{y}) + \int_{C(d) \cap C^*} h(\mathbf{y}) dF(\mathbf{y})$$

and

$$\hat{J}(d^*) = \Gamma_0(1 - p) + \int_{C^* \setminus C(d)} h(\mathbf{y}) dF(\mathbf{y}) + \int_{C(d) \cap C^*} h(\mathbf{y}) dF(\mathbf{y}).$$

Therefore,

$$\hat{J}(d) - \hat{J}(d^*) = \int_{C(d) \setminus C^*} h(\mathbf{y}) dF(\mathbf{y}) + \int_{C^* \setminus C(d)} (-h(\mathbf{y})) dF(\mathbf{y}) \geq 0$$

since

$$\int_{C(d) \setminus C^*} h(\mathbf{y}) dF(\mathbf{y}) \geq 0 \quad \text{and} \quad \int_{C^* \setminus C(d)} h(\mathbf{y}) dF(\mathbf{y}) \leq 0$$

by the very definition of C^* . The problem $\hat{\mathcal{P}}_B$ is therefore solved by the test d^* defined at (2.13). \blacksquare

Uniqueness The solution to the Bayesian problem is *not* unique: It should be plain that C^* could be replaced by

$$C^{**} \equiv \{\mathbf{y} \in \mathbb{R}^k : h(\mathbf{y}) \leq 0\}$$

(with corresponding test d^{**}) without affecting the conclusion of optimality since

$$\int_{\{\mathbf{y} \in \mathbb{R}^k : h(\mathbf{y}) = 0\}} h(\mathbf{y}) dF(\mathbf{y}) = 0.$$

While it is true that $J(d^*) = J(d^{**})$, it is not necessarily the case that the equalities $P_F(d^*) = P_F(d^{**})$ or $P_M(d^*) = P_M(d^{**})$ hold.

Implementation using likelihood ratio test Assume that $0 < p < 1$ to avoid trivial situations, and that the relative costs satisfy the conditions

$$\Gamma_h > 0, \quad h = 0, 1, \tag{2.15}$$

i.e., the cost of making an incorrect decision is greater than the cost of making a correct decision. This is of course a most reasonable assumption which always

holds in applications. Under this condition, the Bayesian decision rule d^* given in Theorem 2.2.1 takes the equivalent form

$$d^*(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < \frac{\Gamma_0(1-p)}{\Gamma_1 p} f_0(\mathbf{y}). \quad (2.16)$$

In view of the definition (1.18), the Bayesian test d^* is indeed a test d_η with η given by

$$\eta \equiv \frac{\Gamma_0(1-p)}{\Gamma_1 p}.$$

Equipped with Lemma 1.6.1 we can now restate Theorem 2.2.1.

Theorem 2.2.2 *Assume the absolute continuity conditions (A.1)–(A.2) to hold. Whenever $\Gamma_h > 0$ for $h = 0, 1$, the Bayesian decision rule d^* identified in Theorem 2.1 is equivalent to the likelihood ratio test Lrt_{η^*} where*

$$\eta^* \equiv \frac{\Gamma_0(1-p)}{\Gamma_1 p} = \frac{C(0,1) - C(0,0)}{C(1,0) - C(1,1)} \cdot \frac{1-p}{p}.$$

2.3 The probability of error criterion

A special case of great interest is obtained when the cost function C takes the form

$$C(h, d) = \mathbf{1}[h \neq d], \quad h, d = 0, 1.$$

The corresponding expected cost then reduces to the probability of making an incorrect decision, namely the *probability of error*, and is given by

$$P_E(d) \equiv \mathbb{P}[d(\mathbf{Y}) \neq H], \quad d \in \mathcal{D}.$$

We check that

$$\Gamma_h = C(h, 1-h) - C(h, h) = 1, \quad h = 0, 1,$$

and the relations (2.5)-(2.6) yield

$$\begin{aligned} P_E(d) &= (1-p) \cdot P_F(d) + p \cdot P_M(d) \\ &= p + (1-p) \cdot P_F(d) - p \cdot P_D(d), \quad d \in \mathcal{D}. \end{aligned} \quad (2.17)$$

For the probability of error criterion, the threshold η^* appearing in Theorem 2.2.2 has the simpler form

$$\eta^* = \frac{1-p}{p}.$$

The optimal decision rule d^* , as described at (2.16), can now be rewritten as

$$d^*(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < \frac{1-p}{p} f_0(\mathbf{y}). \quad (2.18)$$

The ML test In the uniform prior case, i.e., $p = \frac{1}{2}$, the Bayesian test (2.18) becomes

$$d^*(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < f_0(\mathbf{y}). \quad (2.19)$$

In other words, the optimal decision is to select that hypothesis whose likelihood is largest given the observation \mathbf{y} . We refer to this strategy as the *Maximum Likelihood* (ML) test.

The MAP computer Finally, (2.18) can also be rewritten as

$$d^*(\mathbf{y}) = 0 \quad \text{iff} \quad \mathbb{P}[H = 1 | \mathbf{Y} = \mathbf{y}] < \mathbb{P}[H = 0 | \mathbf{Y} = \mathbf{y}] \quad (2.20)$$

since for each \mathbf{y} in \mathbb{R}^k , we have

$$\mathbb{P}[H = 1 | \mathbf{Y} = \mathbf{y}] = \frac{p f_1(\mathbf{y})}{p f_1(\mathbf{y}) + (1-p) f_0(\mathbf{y})}$$

and

$$\mathbb{P}[H = 0 | \mathbf{Y} = \mathbf{y}] = \frac{(1-p) f_0(\mathbf{y})}{p f_1(\mathbf{y}) + (1-p) f_0(\mathbf{y})}$$

by Bayes' Theorem. For each $h = 0, 1$, the conditional probability $\mathbb{P}[H = h | \mathbf{Y} = \mathbf{y}]$ is known as the *posterior* probability that $H = h$ occurs given the observation \mathbf{y} . Put differently, the optimal test (2.20) compares these posterior probabilities given the observation \mathbf{y} , and selects the hypothesis with the largest posterior probability, hence the terminology *Maximum A Posteriori* (MAP) computer.

2.4 The Gaussian case

Assume that the observation rv \mathbf{Y} is conditionally Gaussian given H , i.e.,

$$\begin{aligned} H_1 &: \mathbf{Y} \sim N(\mathbf{m}_1, \mathbf{R}_1) \\ H_0 &: \mathbf{Y} \sim N(\mathbf{m}_0, \mathbf{R}_0) \end{aligned}$$

where \mathbf{m}_1 and \mathbf{m}_0 are elements in \mathbb{R}^k , and the $k \times k$ symmetric matrices \mathbf{R}_1 and \mathbf{R}_0 are *positive definite* (thus *invertible*). Throughout the pairs $(\mathbf{m}_0, \mathbf{R}_0)$ and $(\mathbf{m}_1, \mathbf{R}_1)$ are distinct so that the probability density functions $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ are distinct since

$$f_h(\mathbf{y}) = \frac{1}{\sqrt{(2\pi)^k \det \mathbf{R}_h}} e^{-\frac{1}{2}(\mathbf{y} - \mathbf{m}_h)' \mathbf{R}_h^{-1} (\mathbf{y} - \mathbf{m}_h)}, \quad \mathbf{y} \in \mathbb{R}^k, \quad h = 0, 1.$$

Both conditions **(A.1)** and **(A.2)** obviously hold, and for each $\eta > 0$, the test d_η and Lrt_η coincide.

The likelihood ratio and the likelihood ratio tests For this example, the likelihood ratio function is given by

$$L(\mathbf{y}) = \sqrt{\frac{\det(\mathbf{R}_0)}{\det(\mathbf{R}_1)}} \cdot e^{\frac{1}{2}Q(\mathbf{y})}, \quad \mathbf{y} \in \mathbb{R}^k$$

where we have used the notation

$$Q(\mathbf{y}) = (\mathbf{y} - \mathbf{m}_0)' \mathbf{R}_0^{-1} (\mathbf{y} - \mathbf{m}_0) - (\mathbf{y} - \mathbf{m}_1)' \mathbf{R}_1^{-1} (\mathbf{y} - \mathbf{m}_1).$$

Fix $\eta > 0$. By direct substitution, we conclude that

$$Lrt_\eta(\mathbf{y}) = 0 \quad \text{iff} \quad e^{\frac{1}{2}Q(\mathbf{y})} < \sqrt{\eta^2 \cdot \frac{\det \mathbf{R}_1}{\det \mathbf{R}_0}},$$

and a simple logarithmic transformation yields

$$Lrt_\eta(\mathbf{y}) = 0 \quad \text{iff} \quad Q(\mathbf{y}) < \log \left(\eta^2 \frac{\det \mathbf{R}_1}{\det \mathbf{R}_0} \right).$$

The equal covariance case If the covariances are identical under both hypotheses, i.e.,

$$\mathbf{R}_0 = \mathbf{R}_1 \equiv \mathbf{R},$$

with $\mathbf{m}_1 \neq \mathbf{m}_0$, then

$$\begin{aligned} Q(\mathbf{y}) &= (\mathbf{y} - \mathbf{m}_0)' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{m}_0) - (\mathbf{y} - \mathbf{m}_1)' \mathbf{R}^{-1} (\mathbf{y} - \mathbf{m}_1) \\ &= 2\mathbf{y}' \mathbf{R}^{-1} (\mathbf{m}_1 - \mathbf{m}_0) - (\mathbf{m}_1' \mathbf{R}^{-1} \mathbf{m}_1 - \mathbf{m}_0' \mathbf{R}^{-1} \mathbf{m}_0). \end{aligned} \quad (2.21)$$

The form of Lrt_η simplifies even further to read

$$Lrt_\eta(\mathbf{y}) = 0 \quad \text{iff} \quad \mathbf{y}' \mathbf{R}^{-1} \Delta \mathbf{m} < \tau(\eta)$$

where we have set

$$\Delta \mathbf{m} \equiv \mathbf{m}_1 - \mathbf{m}_0 \quad (2.22)$$

and

$$\tau(\eta) \equiv \frac{1}{2} (\mathbf{m}_1' \mathbf{R}^{-1} \mathbf{m}_1 - \mathbf{m}_0' \mathbf{R}^{-1} \mathbf{m}_0) + \log \eta. \quad (2.23)$$

Evaluating probabilities We will now evaluate the probabilities of false alarm and miss under Lrt_η . It is plain that

$$\begin{aligned} P_F(Lrt_\eta) &= \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 1 | H = 0] \\ &= \mathbb{P}[L(\mathbf{Y}) \geq \eta | H = 0] \\ &= \mathbb{P}[\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} \geq \tau(\eta) | H = 0] \end{aligned} \quad (2.24)$$

and

$$\begin{aligned} P_M(Lrt_\eta) &= \mathbb{P}[Lrt_\eta(\mathbf{Y}) = 0 | H = 1] \\ &= \mathbb{P}[L(\mathbf{Y}) < \eta | H = 1] \\ &= \mathbb{P}[\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} < \tau(\eta) | H = 1] \\ &= 1 - \mathbb{P}[\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} \geq \tau(\eta) | H = 1]. \end{aligned} \quad (2.25)$$

To carry out the calculations further, recall that for each $h = 0, 1$, given $H = h$, the rv \mathbf{Y} is conditionally Gaussian with mean vector \mathbf{m}_h and covariance matrix \mathbf{R} . Therefore, the scalar rv $\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m}$ is also conditionally Gaussian with mean and variance given by

$$\mathbb{E}[\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} | H = h] = \mathbf{m}_h' \mathbf{R}^{-1} \Delta \mathbf{m}$$

and

$$\begin{aligned}
\text{Var} [\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} | H = h] &= (\mathbf{R}^{-1} \Delta \mathbf{m})' \text{Cov} [\mathbf{Y} | H = h] (\mathbf{R}^{-1} \Delta \mathbf{m}) \\
&= (\mathbf{R}^{-1} \Delta \mathbf{m})' \mathbf{R} (\mathbf{R}^{-1} \Delta \mathbf{m}) \\
&= \Delta \mathbf{m}' \mathbf{R}^{-1} \Delta \mathbf{m},
\end{aligned} \tag{2.26}$$

respectively. In obtaining this last relation we have used the fact that

$$\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} = (\mathbf{R}^{-1} \Delta \mathbf{m})' \mathbf{Y}.$$

Consequently, for all $h = 0, 1$,

$$\begin{aligned}
&\mathbb{P} [\mathbf{Y}' \mathbf{R}^{-1} \Delta \mathbf{m} \geq \tau(\eta) | H = h] \\
&= \mathbb{P} \left[\mathbf{m}'_h \mathbf{R}^{-1} \Delta \mathbf{m} + \sqrt{\Delta \mathbf{m}' \mathbf{R}^{-1} \Delta \mathbf{m}} \cdot Z \geq \tau(\eta) \right] \\
&= \mathbb{P} \left[Z \geq \frac{\tau(\eta) - \mathbf{m}'_h \mathbf{R}^{-1} \Delta \mathbf{m}}{\sqrt{\Delta \mathbf{m}' \mathbf{R}^{-1} \Delta \mathbf{m}}} \right]
\end{aligned} \tag{2.27}$$

where $Z \sim N(0, 1)$.

For the sake of convenience, pose

$$d^2 \equiv \Delta \mathbf{m}' \mathbf{R}^{-1} \Delta \mathbf{m}, \tag{2.28}$$

and note that

$$\tau(\eta) - \mathbf{m}'_h \mathbf{R}^{-1} \Delta \mathbf{m} = \begin{cases} \log \eta - \frac{1}{2} d^2 & \text{if } h = 1 \\ \log \eta + \frac{1}{2} d^2 & \text{if } h = 0. \end{cases}$$

It is now clear that

$$P_F(Lrt_\eta) = 1 - \Phi \left(\frac{\log \eta + \frac{1}{2} d^2}{d} \right)$$

and

$$P_M(Lrt_\eta) = \Phi \left(\frac{\log \eta - \frac{1}{2} d^2}{d} \right).$$

We finally obtain

$$P_D(Lrt_\eta) = 1 - \Phi \left(\frac{\log \eta - \frac{1}{2} d^2}{d} \right).$$

The ML test The ML test corresponds to $\eta = 1$, in which case these expressions become

$$P_F(d_{\text{ML}}) = 1 - \Phi\left(\frac{d}{2}\right) = Q\left(\frac{d}{2}\right)$$

and

$$P_M(d_{\text{ML}}) = \Phi\left(-\frac{d}{2}\right) = Q\left(\frac{d}{2}\right),$$

whence

$$P_E(d_{\text{ML}}) = (1 - p)P_F(d_{\text{ML}}) + pP_M(d_{\text{ML}}) = Q\left(\frac{d}{2}\right)$$

regardless of the prior p .

2.5 The Bernoulli case

Consider the binary hypothesis testing problem

$$\begin{aligned} H_1 : Y &\sim \text{Ber}(a_1) \\ H_0 : Y &\sim \text{Ber}(a_0) \end{aligned}$$

with $a_1 < a_0$ in $(0, 1)$. The case $a_0 < a_1$ is left as an exercise. Thus,

$$\mathbb{P}[Y = 1|H = h] = a_h = 1 - \mathbb{P}[Y = 0|H = h], \quad h = 0, 1$$

and Conditions **(A.1)** and **(A.2)** obviously hold with respect to counting measure F on $\{0, 1\}$. The likelihood rate function is given by

$$L(y) = \left(\frac{1 - a_1}{1 - a_0}\right)^{1-y} \left(\frac{a_1}{a_0}\right)^y, \quad y \in \mathbb{R}.$$

For each $\eta > 0$, the test d_η takes the following form

$$\begin{aligned} d_\eta(y) = 0 & \quad \text{iff} \quad \left(\frac{1 - a_1}{1 - a_0}\right)^{1-y} \cdot \left(\frac{a_1}{a_0}\right)^y < \eta \\ & \quad \text{iff} \quad \left(\frac{1 - a_0}{1 - a_1} \cdot \frac{a_1}{a_0}\right)^y < \eta \cdot \frac{1 - a_0}{1 - a_1}, \quad y \in \mathbb{R}. \end{aligned} \quad (2.29)$$

Therefore,

$$\begin{aligned}
P_F(d_\eta) &= \mathbb{P}[d_\eta(Y) = 1|H = 0] \\
&= \mathbb{P}\left[\left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y \geq \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 0\right] \\
&= \mathbb{P}\left[Y = 1, \left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y \geq \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 0\right] \\
&\quad + \mathbb{P}\left[Y = 0, \left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y \geq \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 0\right] \\
&= a_0 \mathbf{1}\left[\eta \cdot \frac{1-a_0}{1-a_1} \leq \frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right] + (1-a_0) \mathbf{1}\left[\eta \cdot \frac{1-a_0}{1-a_1} \leq 1\right] \\
&= a_0 \mathbf{1}\left[\eta \leq \frac{a_1}{a_0}\right] + (1-a_0) \mathbf{1}\left[\eta \frac{1-a_0}{1-a_1} \leq 1\right] \\
&= a_0 \mathbf{1}\left[\eta \leq \frac{a_1}{a_0}\right] + (1-a_0) \mathbf{1}\left[\eta \leq \frac{1-a_1}{1-a_0}\right]. \tag{2.30}
\end{aligned}$$

Similarly, we get

$$\begin{aligned}
P_M(d_\eta) &= \mathbb{P}[d_\eta(Y) = 0|H = 1] \\
&= \mathbb{P}\left[\left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y < \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 1\right] \\
&= \mathbb{P}\left[Y = 1, \left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y < \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 1\right] \\
&\quad + \mathbb{P}\left[Y = 0, \left(\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0}\right)^Y < \eta \cdot \frac{1-a_0}{1-a_1} \middle| H = 1\right] \\
&= a_1 \mathbf{1}\left[\frac{1-a_0}{1-a_1} \cdot \frac{a_1}{a_0} < \eta \cdot \frac{1-a_0}{1-a_1}\right] + (1-a_1) \mathbf{1}\left[1 < \eta \cdot \frac{1-a_0}{1-a_1}\right] \\
&= a_1 \mathbf{1}\left[\frac{a_1}{a_0} < \eta\right] + (1-a_1) \mathbf{1}\left[1 < \eta \cdot \frac{1-a_0}{1-a_1}\right] \\
&= a_1 \mathbf{1}\left[\frac{a_1}{a_0} < \eta\right] + (1-a_1) \mathbf{1}\left[\frac{1-a_1}{1-a_0} < \eta\right]. \tag{2.31}
\end{aligned}$$

2.6 Additional examples

We now present several examples where Conditions **(A.1)** or **(A.2)** fail. In all cases we assume $\Gamma_0 > 0$ and $\Gamma_1 > 0$.

An example where absolute continuity (A.2) fails Here, the observation is the scalar rv Y with F_0 and F_1 admitting probability density functions $f_0, f_1 : \mathbb{R} \rightarrow \mathbb{R}_+$ with respect to Lebesgue measure given by

$$f_0(y) = \begin{cases} 1 - |y| & \text{if } |y| \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad f_1(y) = \begin{cases} \frac{1}{3} & \text{if } -1 \leq y \leq 2 \\ 0 & \text{otherwise.} \end{cases}$$

Condition **(A.1)** holds (with Lebesgue measure) but the absolute continuity condition **(A.2)** is clearly not satisfied. However, simple substitution reveals that

$$\begin{aligned} h(y) &= \Gamma_1 p \cdot f_1(y) - \Gamma_0(1-p) \cdot f_0(y) \\ &= \begin{cases} 0 & \text{if } y < -1 \\ \frac{1}{3}\Gamma_1 p - \Gamma_0(1-p)(1-|y|) & \text{if } |y| \leq 1 \\ \frac{1}{3}\Gamma_1 p & \text{if } 1 < y \leq 2 \\ 0 & \text{if } 2 < y. \end{cases} \end{aligned} \quad (2.32)$$

The Bayesian test d^* is simply

$$d^*(y) = 0 \quad \text{iff} \quad |y| < 1 - \frac{\frac{1}{3}\Gamma_1 p}{\Gamma_0(1-p)}.$$

Another example where absolute continuity (A.2) fails The observation is the scalar rv Y with F_0 and F_1 admitting probability density functions $f_0, f_1 : \mathbb{R} \rightarrow \mathbb{R}_+$ with respect to Lebesgue measure given by

$$f_0(y) = \begin{cases} 1 - |y| & \text{if } |y| \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad f_1(y) = \begin{cases} \frac{1}{3} & \text{if } 0 \leq y \leq 3 \\ 0 & \text{otherwise.} \end{cases}$$

Condition **(A.1)** holds (with Lebesgue measure) but **(A.2)** fails. Simple substitution reveals that

$$\begin{aligned}
 h(y) &= \Gamma_1 p \cdot f_1(y) - \Gamma_0(1-p) \cdot f_0(y) \\
 &= \begin{cases} 0 & \text{if } y < -1 \\ -\Gamma_0(1-p)(1+y) & \text{if } -1 \leq y \leq 0 \\ \frac{1}{3}\Gamma_1 p - \Gamma_0(1-p)(1-y) & \text{if } 0 < y \leq 1 \\ \frac{1}{3}\Gamma_1 p & \text{if } 1 < y \leq 3 \\ 0 & \text{if } 3 < y, \end{cases} \quad (2.33)
 \end{aligned}$$

and it is straightforward to check that the Bayesian test d^* is simply

$$d^*(y) = 0 \quad \text{iff} \quad \begin{array}{c} -1 < y \leq 0 \\ \text{or} \\ 0 < y \leq 1, y < 1 - \frac{\frac{1}{3}\Gamma_1 p}{\Gamma_0(1-p)}. \end{array}$$

Equivalently, d^* can be described as

$$d^*(y) = 0 \quad \text{iff} \quad y \in \left(-1, \left(1 - \frac{\Gamma_1 p}{3\Gamma_0(1-p)} \right)^+ \right).$$

A final example Consider the binary hypothesis testing problem

$$\begin{aligned}
 H_1 : & Y \sim F_1 \\
 H_0 : & Y \sim F_0
 \end{aligned}$$

where F_0 is the discrete uniform distribution on $\{0, 1\}$, and F_1 is uniform on the interval $(0, 1)$. Thus, F_1 admits a probability density function $f_1 : \mathbb{R} \rightarrow \mathbb{R}_+$ with respect to Lebesgue measure given by

$$f_1(y) = \begin{cases} 1 & \text{if } y \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

and

$$\mathbb{P}[Y = 0 | H = 0] = \mathbb{P}[Y = 1 | H = 0] = \frac{1}{2}.$$

In this example F cannot be taken to be either the distribution associated with Lebesgue measure on \mathbb{R} or with the counting measure on $\{0, 1\}$. In principle we could use F given by (1.10) but this would yield complicated expressions for the density functions $f_0, f_1 : \mathbb{R} \rightarrow \mathbb{R}^+$. Instead of applying Theorem 2.2.1 with that choice, we provide a direct optimization of the auxiliary expected cost function (2.4): For each test d in \mathcal{D} we recall that we have

$$\begin{aligned} & \widehat{J}(d) \\ = & \Gamma_0(1-p) + \Gamma_1 p \cdot \mathbb{P}[d(Y) = 0 | H = 1] - \Gamma_0(1-p) \cdot \mathbb{P}[d(Y) = 0 | H = 0] \end{aligned}$$

with

$$\mathbb{P}[d(Y) = 0 | H = 0] = \begin{cases} \frac{1}{2} & \text{if } 0 \in C(d), 1 \notin C(d) \\ \frac{1}{2} & \text{if } 1 \in C(d), 0 \notin C(d) \\ 1 & \text{if } 0 \in C(d), 1 \in C(d) \\ 0 & \text{if } 0 \notin C(d), 1 \notin C(d) \end{cases}$$

and

$$\mathbb{P}[d(Y) = 0 | H = 1] = \int_{C(d)} f_1(y) dy = |C(d) \cap [0, 1]|.$$

Adding or deleting a finite number of points from $C(d)$ will *not* affect the value of $\mathbb{P}[d(Y) = 0 | H = 1]$, but it may change the value of $\mathbb{P}[d(Y) = 0 | H = 0]$. Therefore, with $C(d)$ given, modify it, if needed, by adding both points 0 and 1. If C' denotes this Borel subset of \mathbb{R} , then $C' = C(d) \cup \{0, 1\}$; if d' denotes the corresponding test, then $C(d') = C'$. Obviously

$$\mathbb{P}[d(Y) = 0 | H = 1] = \mathbb{P}[d'(Y) = 0 | H = 1] = |C(d') \cap [0, 1]|$$

since $|C(d') \cap [0, 1]| = |C(d) \cap [0, 1]|$, while

$$\mathbb{P}[d(Y) = 0 | H = 0] \leq \mathbb{P}[d'(Y) = 0 | H = 0] = 1.$$

We can now conclude that

$$\begin{aligned} & \widehat{J}(d) \\ = & \Gamma_0(1-p) + \Gamma_1 p \cdot \mathbb{P}[d(Y) = 0 | H = 1] - \Gamma_0(1-p) \cdot \mathbb{P}[d(Y) = 0 | H = 0] \\ \geq & \Gamma_0(1-p) + \Gamma_1 p \cdot \mathbb{P}[d'(Y) = 0 | H = 1] - \Gamma_0(1-p) \cdot \mathbb{P}[d'(Y) = 0 | H = 0] \\ = & \Gamma_0(1-p) + \Gamma_1 p \cdot |C(d') \cap [0, 1]| - \Gamma_0(1-p) \\ = & \Gamma_1 p \cdot |C(d') \cap [0, 1]| \geq 0. \end{aligned} \tag{2.34}$$

Consider the test $d^* : \mathbb{R} \rightarrow \{0, 1\}$ given by

$$d^*(y) = \mathbf{1}_{[\{0, 1\}]}(y), \quad y \in \mathbb{R}.$$

The arguments leading to (2.34) also show that

$$\hat{J}(d^*) = \Gamma_1 p \cdot |C(d^*) \cap [0, 1]| = 0,$$

and the test d^* is therefore a Bayesian decision rule.

2.7 Exercises

2.8 References

Chapter 3

Randomized tests

As we shall see shortly, a solution cannot always be found to the Minimax and Neyman–Pearson formulations of the hypothesis testing problem if the search is restricted to the class of decision rules \mathcal{D} as done for the Bayesian set–up. In some very real sense this class \mathcal{D} of tests is not always large enough to guarantee a solution; to remedy this difficulty we enlarge \mathcal{D} by considering the class of *randomized* tests or decision rules.

3.1 Randomized tests

We start with a definition.

A *randomized* test δ is a Borel mapping $\delta : \mathbb{R}^k \rightarrow [0, 1]$ with the following interpretation as conditional probability: Having observed $\mathbf{Y} = \mathbf{y}$, it is decided that the state of nature is 1 (resp. 0) with probability $\delta(\mathbf{y})$ (resp. $1 - \delta(\mathbf{y})$). The collection of all randomized tests will be denoted by \mathcal{D}^* .

Obviously, any test d in \mathcal{D} can be mechanized as a randomized test, say $\delta_d : \mathbb{R}^k \rightarrow [0, 1]$, given by

$$\delta_d(\mathbf{y}) \equiv d(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k.$$

A test in \mathcal{D} is often referred to as a *pure* strategy.

A natural question then arises as to how such randomization mechanisms can be incorporated into the probabilistic framework introduced earlier in Section 1.2: The model data is unchanged as we are given two probability distributions F_0 and F_1 on \mathbb{R}^k and a prior p in $[0, 1]$. We still consider a sample space Ω equipped with

a σ -field of events \mathcal{F} , and on it we now define the three rvs H , \mathbf{Y} and D which take values in $\{0, 1\}$, \mathbb{R}^k and $\{0, 1\}$, respectively. The rvs H and \mathbf{Y} have the same interpretation as before, as state of nature and observation, respectively, while the rv D now encodes the decision to be taken on the basis of the observation \mathbf{Y} .

With each decision rule δ in \mathcal{D}^* we associate a probability measure \mathbb{P}_δ on \mathcal{F} such that the following constraints are satisfied: As before, this time under \mathbb{P}_δ , we still have

$$\mathbb{P}_\delta[\mathbf{Y} \leq \mathbf{y} | H = h] = F_h(\mathbf{y}), \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ h = 0, 1 \end{array}$$

and

$$p = \mathbb{P}_\delta[H = 1] = 1 - \mathbb{P}_\delta[H = 0].$$

Therefore, under \mathbb{P}_δ the probability distribution of the pair (H, \mathbf{Y}) does not depend on δ with

$$\mathbb{P}_\delta[H = h, \mathbf{Y} \leq \mathbf{y}] = \mathbb{P}_\delta[H = h] F_h(\mathbf{y}), \quad \begin{array}{l} h = 0, 1, \\ \mathbf{y} \in \mathbb{R}^k \end{array} \quad (3.1)$$

as expected. In addition, for $h = 0, 1$ and \mathbf{y} in \mathbb{R}^k , we now require that

$$\begin{aligned} \mathbb{P}_\delta[D = d | H = h, \mathbf{Y} = \mathbf{y}] &= \begin{cases} 1 - \delta(\mathbf{y}) & \text{if } d = 0 \\ \delta(\mathbf{y}) & \text{if } d = 1 \end{cases} \\ &= d\delta(\mathbf{y}) + (1 - d)(1 - \delta(\mathbf{y})). \end{aligned} \quad (3.2)$$

The joint probability distribution of the rvs H , D and \mathbf{Y} (under \mathbb{P}_δ) can now be completely specified: With $h, d = 0, 1$ and a Borel subset B of \mathbb{R}^k , a preconditioning argument gives

$$\begin{aligned} &\mathbb{P}_\delta[H = h, D = d, \mathbf{Y} \in B] \\ &= \mathbb{E}_\delta[\mathbf{1}[H = h, \mathbf{Y} \in B] \mathbb{P}_\delta[D = d | H, \mathbf{Y}]] \\ &= \mathbb{E}_\delta[\mathbf{1}[H = h, \mathbf{Y} \in B] (d\delta(\mathbf{Y}) + (1 - d)(1 - \delta(\mathbf{Y})))] \\ &= \mathbb{P}_\delta[H = h] \cdot \int_B (d\delta(\mathbf{y}) + (1 - d)(1 - \delta(\mathbf{y}))) dF_h(\mathbf{y}) \\ &= \begin{cases} \mathbb{P}_\delta[H = h] \cdot \int_B (1 - \delta(\mathbf{y})) dF_h(\mathbf{y}) & \text{if } d = 0 \\ \mathbb{P}_\delta[H = h] \cdot \int_B \delta(\mathbf{y}) dF_h(\mathbf{y}) & \text{if } d = 1. \end{cases} \end{aligned} \quad (3.3)$$

3.2 An alternate framework

The class \mathcal{D}^* of randomized strategies gives rise to a *collection* of probability triples, namely

$$\{(\Omega, \mathcal{F}, \mathbb{P}_\delta), \delta \in \mathcal{D}^*\}.$$

It is however possible to provide an equivalent probabilistic framework using a *single* probability triple $(\Omega, \mathcal{F}, \mathbb{P})$. To see how this can be done, imagine that the original probability triple $(\Omega, \mathcal{F}, \mathbb{P})$ is sufficiently rich that there exists on it a rv $U : \Omega \rightarrow [0, 1]$ which is uniformly distributed on $(0, 1)$, and independent of the pair of rvs H and \mathbf{Y} . This amounts to

$$\mathbb{P}[U \leq t, H = h, \mathbf{Y} \leq \mathbf{y}] = \mathbb{P}[U \leq t] \mathbb{P}[H = h, \mathbf{Y} \leq \mathbf{y}], \quad \begin{array}{l} t \in \mathbb{R} \\ h = 0, 1, \\ \mathbf{y} \in \mathbb{R}^k \end{array}$$

with

$$\mathbb{P}[U \leq t] = \begin{cases} 0 & \text{if } t \leq 0 \\ \min(t, 1) & \text{if } t \geq 0, \end{cases}$$

$$\mathbb{P}[H = h, \mathbf{Y} \leq \mathbf{y}] = \mathbb{P}[H = h] F_h(\mathbf{y}), \quad \begin{array}{l} h = 0, 1, \\ \mathbf{y} \in \mathbb{R}^k \end{array}$$

and

$$\mathbb{P}[H = 1] = p = 1 - \mathbb{P}[H = 0].$$

Now, for each decision rule δ in \mathcal{D}^* , define the $\{0, 1\}$ -valued rv D_δ given by

$$D_\delta = \mathbf{1}[U \leq \delta(\mathbf{Y})].$$

Note that

$$\begin{aligned} \mathbb{P}[D_\delta = 1 | H = h, \mathbf{Y} = \mathbf{y}] &= \mathbb{E}[\mathbf{1}[U \leq \delta(\mathbf{Y})] | H = h, \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{E}[\mathbf{1}[U \leq \delta(\mathbf{y})] | H = h, \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{P}[U \leq \delta(\mathbf{y})] \\ &= \delta(\mathbf{y}) \end{aligned} \tag{3.4}$$

under the enforced independence assumptions. Similarly it follows that

$$\begin{aligned} \mathbb{P}[D_\delta = 0 | H = h, \mathbf{Y} = \mathbf{y}] &= 1 - \mathbb{P}[D_\delta = 1 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= 1 - \delta(\mathbf{y}). \end{aligned} \tag{3.5}$$

Therefore, the conditional distribution of D_δ (under \mathbb{P}) given H and \mathbf{Y} coincides with the conditional distribution of D (under \mathbb{P}_δ) given H and \mathbf{Y} , and the two formalisms are probabilistically equivalent.

3.3 Evaluating error probabilities

Evaluating error probabilities under randomized tests can be done rather easily: Consider a randomized test δ in \mathcal{D}^* . In analogy with (1.15) and (1.16), we evaluate the probabilities of false alarm and miss under δ as

$$P_F(\delta) \equiv \mathbb{P}_\delta [D = 1 | H = 0] \quad (3.6)$$

and

$$P_M(\delta) \equiv \mathbb{P}_\delta [D = 0 | H = 1]. \quad (3.7)$$

It is also convenient to consider the so-called probability of *detection* given by

$$P_D(\delta) \equiv \mathbb{P}_\delta [D = 1 | H = 1] = 1 - P_M(\delta). \quad (3.8)$$

Because

$$\mathbb{P}_\delta [D = h | H] = \mathbb{E}_\delta [\mathbb{P}_\delta [D = h | H, \mathbf{Y}] | H], \quad h = 0, 1$$

we readily conclude that

$$P_F(\delta) = \int_{\mathbb{R}^k} \delta(\mathbf{y}) dF_0(\mathbf{y}) \quad (3.9)$$

and

$$P_M(\delta) = \int_{\mathbb{R}^k} (1 - \delta(\mathbf{y})) dF_1(\mathbf{y}), \quad (3.10)$$

so that

$$P_D(\delta) = \int_{\mathbb{R}^k} \delta(\mathbf{y}) dF_1(\mathbf{y}), \quad (3.11)$$

3.4 The Bayesian problem revisited

Assuming the cost function $C : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$ introduced in Section 2.1, we define the expected cost function $J^* : \mathcal{D}^* \rightarrow \mathbb{R}$ given by

$$J^*(\delta) = \mathbb{E}_\delta [C(H, D)], \quad \delta \in \mathcal{D}^*.$$

When considering randomized decision rules, the original Bayesian Problem \mathcal{P}_B is now reformulated as the minimization problem

$$\mathcal{P}_B^* : \quad \text{Minimize } J^*(\delta) \text{ over } \delta \text{ in } \mathcal{D}^*.$$

This amounts to finding an admissible test δ^* in \mathcal{D}^* such that

$$J^*(\delta^*) \leq J^*(\delta), \quad \delta \in \mathcal{D}^*. \quad (3.12)$$

Any admissible test δ^* which satisfies (3.12) is called a randomized Bayesian test, and the value

$$J^*(\delta^*) = \inf_{\delta \in \mathcal{D}^*} J^*(\delta) \quad (3.13)$$

is sometimes referred to as the randomized *Bayesian cost*.

Obviously, since $\mathcal{D} \subset \mathcal{D}^*$ (with a slight abuse of notation) with

$$J^*(\delta_d) = J(d), \quad d \in \mathcal{D},$$

it is plain that

$$\inf_{\delta \in \mathcal{D}^*} J^*(\delta) \leq \inf_{d \in \mathcal{D}} J(d).$$

While in principle this last inequality could be strict, we now show that it is not so and that the Bayesian problem is not affected by considering the larger set of randomized decision rules; the proof is available in Section 3.6.

Theorem 3.4.1 *Under the absolute continuity condition (A.1), it holds that*

$$\inf_{\delta \in \mathcal{D}^*} J^*(\delta) = \inf_{d \in \mathcal{D}} J(d). \quad (3.14)$$

It follows from Theorem 2.13 that (3.14) is equivalent to

$$\min_{\delta \in \mathcal{D}^*} J^*(\delta) = \min_{d \in \mathcal{D}} J(d) = J(d^*) \quad (3.15)$$

where the deterministic test $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ is given by (2.13).

For easy reference we close with the following analog of Lemma 2.1.1 for randomized tests; the proof is left as an exercise.

Lemma 3.4.1 For any admissible rule δ in \mathcal{D}^* , the relation

$$J^*(\delta) = \mathbb{E}[C(H, H)] + \widehat{J}^*(\delta) \quad (3.16)$$

holds with

$$\widehat{J}^*(\delta) = \Gamma_0(1 - p) \cdot P_F(\delta) + \Gamma_1 p \cdot P_M(\delta). \quad (3.17)$$

3.5 Randomizing between two pure decision rules

Consider two pure strategies d_1 and d_2 in \mathcal{D} . With a in $(0, 1)$, we introduce a randomized policy δ_a in \mathcal{D}^* which first selects the pure strategy d_1 (resp. d_2) with probability a (resp. $1 - a$), and then uses the pure policy that was selected. Formally, this amounts to defining $\delta_a : \mathbb{R}^k \rightarrow [0, 1]$ by

$$\delta_a(\mathbf{y}) = ad_1(\mathbf{y}) + (1 - a)d_2(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k.$$

Applying the expressions (3.9) and (3.10) with the randomized test δ_a we get

$$\begin{aligned} P_F(\delta_a) &= \int_{\mathbb{R}^k} \delta_a(\mathbf{y}) dF_0(\mathbf{y}) \\ &= \int_{\mathbb{R}^k} (ad_1(\mathbf{y}) + (1 - a)d_2(\mathbf{y})) dF_0(\mathbf{y}) \\ &= a \int_{\mathbb{R}^k} d_1(\mathbf{y}) dF_0(\mathbf{y}) + (1 - a) \int_{\mathbb{R}^k} d_2(\mathbf{y}) dF_0(\mathbf{y}) \\ &= aP_F(d_1) + (1 - a)P_F(d_2). \end{aligned} \quad (3.18)$$

Similarly we find that

$$\begin{aligned} P_M(\delta_a) &= \int_{\mathbb{R}^k} (1 - \delta_a(\mathbf{y})) dF_1(\mathbf{y}) \\ &= \int_{\mathbb{R}^k} (1 - ad_1(\mathbf{y}) - (1 - a)d_2(\mathbf{y})) dF_1(\mathbf{y}) \\ &= a \int_{\mathbb{R}^k} (1 - d_1(\mathbf{y})) dF_1(\mathbf{y}) + (1 - a) \int_{\mathbb{R}^k} (1 - d_2(\mathbf{y})) dF_1(\mathbf{y}) \\ &= aP_M(d_1) + (1 - a)P_M(d_2). \end{aligned} \quad (3.19)$$

It immediately follows from (3.16) and (3.17) that

$$J_p^*(\delta_a) = aJ_p(d_1) + (1 - a)J_p(d_2). \quad (3.20)$$

as we use the relations (3.18) and (3.19).

One very concrete way to implement the randomized policy δ_a on the original triple $(\Omega, \mathcal{F}, \mathbb{P})$ proceeds as follows: Consider the original probabilistic framework introduced in Section 1.2 and assume it to be sufficiently rich to carry an additional \mathbb{R} -valued rv V which is independent of the rvs H and \mathbf{Y} (under \mathbb{P}), and is uniformly distributed on the interval $[0, 1]$. Define the $\{0, 1\}$ -valued rv B_a given by

$$B_a = \mathbf{1}[V \leq a].$$

It is plain that the rv B_a is independent of the rvs H and \mathbf{Y} (under \mathbb{P}), with

$$\mathbb{P}[B_a = 1] = a = 1 - \mathbb{P}[B_a = 0].$$

Define the decision rv D_a given by

$$D_a = B_a d_1(\mathbf{Y}) + (1 - B_a) d_2(\mathbf{Y}).$$

It is easy to check that

$$\begin{aligned} & \mathbb{P}[D_a = 1 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{P}[B_a d_1(\mathbf{Y}) + (1 - B_a) d_2(\mathbf{Y}) = 1 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{P}[B_a d_1(\mathbf{y}) + (1 - B_a) d_2(\mathbf{y}) = 1 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= \mathbb{P}[B_a = 1, d_1(\mathbf{y}) = 1 | H = h, \mathbf{Y} = \mathbf{y}] + \mathbb{P}[B_a = 0, d_2(\mathbf{y}) = 1 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= d_1(\mathbf{y}) \mathbb{P}[B_a = 1 | H = h, \mathbf{Y} = \mathbf{y}] + d_2(\mathbf{y}) \mathbb{P}[B_a = 0 | H = h, \mathbf{Y} = \mathbf{y}] \\ &= d_1(\mathbf{y}) \mathbb{P}[B_a = 1] + d_2(\mathbf{y}) \mathbb{P}[B_a = 0] \\ &= ad_1(\mathbf{y}) + (1 - a)d_2(\mathbf{y}), \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ h = 0, 1 \end{array} \end{aligned} \quad (3.21)$$

as desired.

3.6 A proof of Theorem 3.4.1

Pick an arbitrary test δ in \mathcal{D}^* . A simple preconditioning argument shows that

$$J^*(\delta) = \mathbb{E}_\delta[C(H, D)]$$

$$\begin{aligned}
&= \mathbb{E}_\delta [\mathbb{E}_\delta [C(H, D)|H, \mathbf{Y}]] \\
&= \mathbb{E}_\delta [C(H, 1)\mathbb{P}_\delta [D = 1|H, \mathbf{Y}] + C(H, 0)\mathbb{P}_\delta [D = 0|H, \mathbf{Y}]] \\
&= \mathbb{E}_\delta [C(H, 1) \cdot \delta(\mathbf{Y}) + C(H, 0) \cdot (1 - \delta(\mathbf{Y}))] \\
&= \mathbb{E}_\delta [C(H, 0)] + \mathbb{E}_\delta [(C(H, 1) - C(H, 0)) \cdot \delta(\mathbf{Y})] \tag{3.22}
\end{aligned}$$

with

$$\begin{aligned}
&\mathbb{E}_\delta [(C(H, 1) - C(H, 0)) \cdot \delta(\mathbf{Y})] \\
&= \mathbb{E}_\delta [(C(H, 1) - C(H, 0)) \cdot \mathbb{E}_\delta [\delta(\mathbf{Y})|H]] \\
&= (C(1, 1) - C(1, 0)) \mathbb{E}_\delta [\delta(\mathbf{Y})|H = 1] \mathbb{P}_\delta [H = 1] \\
&\quad + (C(0, 1) - C(0, 0)) \mathbb{E}_\delta [\delta(\mathbf{Y})|H = 0] \mathbb{P}_\delta [H = 0] \\
&= -\Gamma_1 p \cdot \mathbb{E}_\delta [\delta(\mathbf{Y})|H = 1] + \Gamma_0(1 - p) \cdot \mathbb{E}_\delta [\delta(\mathbf{Y})|H = 0]. \tag{3.23}
\end{aligned}$$

Using the absolute continuity condition **(A.1)** we can now write

$$\mathbb{E}_\delta [\delta(\mathbf{Y})|H = h] = \int_{\mathbb{R}^k} \delta(\mathbf{y}) dF_h(\mathbf{y}) = \int_{\mathbb{R}^k} \delta(\mathbf{y}) f_h(\mathbf{y}) dF(\mathbf{y}), \quad h = 0, 1$$

so that

$$\begin{aligned}
&J^*(\delta) - \mathbb{E}_\delta [C(H, 0)] \\
&= -\Gamma_1 p \cdot \int_{\mathbb{R}^k} \delta(\mathbf{y}) f_1(\mathbf{y}) dF(\mathbf{y}) + \Gamma_0(1 - p) \cdot \int_{\mathbb{R}^k} \delta(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) \\
&= \int_{\mathbb{R}^k} (-\Gamma_1 p f_1(\mathbf{y}) + \Gamma_0(1 - p) f_0(\mathbf{y})) \delta(\mathbf{y}) dF(\mathbf{y}) \\
&= - \int_{\mathbb{R}^k} h(\mathbf{y}) \delta(\mathbf{y}) dF(\mathbf{y}) \tag{3.24}
\end{aligned}$$

where the mapping $h : \mathbb{R}^k \rightarrow \mathbb{R}$ is given by (2.11). Note that the term $\mathbb{E}_\delta [C(H, 0)]$ does not depend on the randomized test δ being used.

From Theorem 2.2.1 recall that the Bayesian rule which solves Problem \mathcal{P}_B is the test $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ in \mathcal{D} given by (2.13). Note that d^* can also be interpreted as the randomized rule $\delta^* : \mathbb{R}^k \rightarrow [0, 1]$ given by

$$\delta^*(\mathbf{y}) = \begin{cases} 0 & \text{if } h(\mathbf{y}) < 0 \\ 1 & \text{if } h(\mathbf{y}) \geq 0 \end{cases} = \begin{cases} 0 & \text{if } \mathbf{y} \in C^* \\ 1 & \text{if } \mathbf{y} \notin C^* \end{cases}$$

where C^* is the Borel subset of \mathbb{R}^k given by (2.12).

The desired result will be established if we show that

$$J^*(\delta^*) \leq J^*(\delta), \quad \delta \in \mathcal{D}^*.$$

The approach we take is reminiscent of the one used in the proof of Theorem 2.2.1: For an arbitrary δ in \mathcal{D}^* , earlier calculations (3.24) show that

$$\begin{aligned} J^*(\delta) - J^*(\delta^*) &= - \int_{\mathbb{R}^k} h(\mathbf{y})\delta(\mathbf{y})dF(\mathbf{y}) + \int_{\mathbb{R}^k} h(\mathbf{y})\delta^*(\mathbf{y})dF(\mathbf{y}) \\ &= \int_{\mathbb{R}^k} h(\mathbf{y}) (\delta^*(\mathbf{y}) - \delta(\mathbf{y})) dF(\mathbf{y}) \\ &= \int_{C^*} (-h(\mathbf{y}))\delta(\mathbf{y})dF(\mathbf{y}) + \int_{\mathbb{R}^k \setminus C^*} (1 - \delta(\mathbf{y})) h(\mathbf{y})dF(\mathbf{y}) \\ &\geq 0 \end{aligned}$$

as desired since

$$\int_{C^*} (-h(\mathbf{y}))dF(\mathbf{y}) \geq 0 \quad \text{and} \quad \int_{\mathbb{R}^k \setminus C^*} (1 - \delta(\mathbf{y})) h(\mathbf{y})dF(\mathbf{y}) \geq 0$$

by the very definition of the set C^* and of the mapping $h : \mathbb{R}^k \rightarrow \mathbb{R}$. ■

3.7 Exercises

3.8 References

Chapter 4

The Minimax formulation

The Bayesian formulation *implicitly* assumes knowledge of the prior distribution on the hypothesis rv H . In many situations, this assumption cannot be adequately justified, and the Bayesian formulation has to be abandoned for the so-called *Minimax formulation* discussed in this chapter.

4.1 Keeping track of the prior

To facilitate the discussion, we augment the notation introduced in Chapter 1 and Chapter 3 by explicitly indicating the dependence on the prior probability distribution: As before we are given two distinct probability distributions $F_0, F_1 : \mathbb{R}^k \rightarrow [0, 1]$ which act as conditional probability distributions for the observation given the state of nature. As in Chapter 1, we can always construct a collection $\{(\Omega, \mathcal{F}, \mathbb{P}_p), p \in [0, 1]\}$ of probability triples, and rvs H and \mathbf{Y} defined on Ω which take values in $\{0, 1\}$ and \mathbb{R}^k , respectively, such that for each p in $[0, 1]$,

$$F_h(\mathbf{y}) = \mathbb{P}_p[\mathbf{Y} \leq \mathbf{y} | H = h], \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ h = 0, 1 \end{array}$$

and

$$p = \mathbb{P}_p[H = 1] = 1 - \mathbb{P}_p[H = 0].$$

One possible construction was given in Section 1.3: Take $\Omega = \{0, 1\} \times \mathbb{R}^k$ with generic element $\omega = (h, \mathbf{y})$ with $h = 0, 1$ and \mathbf{y} an arbitrary element of \mathbb{R}^k . We endow Ω with the σ -field \mathcal{F} given by

$$\mathcal{F} = \sigma(\mathcal{P}(\{0, 1\}) \times \mathcal{B}(\mathbb{R}^k))$$

where $\mathcal{P}(\{0, 1\})$ is the power set of $\{0, 1\}$, and $\mathcal{B}(\mathbb{R}^k)$ is the Borel σ -field on \mathbb{R}^k . We define the mappings $H : \Omega \rightarrow \mathbb{R}$ and $\mathbf{Y} : \Omega \rightarrow \mathbb{R}^k$ by

$$H(\omega) = h \quad \text{and} \quad \mathbf{Y}(\omega) = \mathbf{y}, \quad \omega = (h, \mathbf{y}) \in \Omega.$$

Both projection mappings are Borel measurable, and therefore define rvs.

As before, it is plain that

$$\begin{aligned} \mathbb{P}_p[\mathbf{Y} \leq \mathbf{y}, H = h] &= \mathbb{P}_p[\mathbf{Y} \leq \mathbf{y} | H = h] \mathbb{P}_p[H = h] \\ &= \begin{cases} (1-p)F_0(\mathbf{y}) & \text{if } h = 0, \mathbf{y} \in \mathbb{R}^k \\ pF_1(\mathbf{y}) & \text{if } h = 1, \mathbf{y} \in \mathbb{R}^k. \end{cases} \end{aligned} \quad (4.1)$$

Let $\mathbb{E}_p[\cdot]$ denote expectation under \mathbb{P}_p .

When dealing with randomized strategies we further augment the notation \mathbb{P}_δ to read $\mathbb{P}_{\delta,p}$ when using the randomized strategy δ in \mathcal{D}^* with prior p ; see Section 3.1 for details on the probabilistic framework to be used.. In that case let $\mathbb{E}_{\delta,p}[\cdot]$ denote expectation under $\mathbb{P}_{\delta,p}$.

4.2 The Bayesian problems

Fix p in $[0, 1]$. Let $J_p(d)$ denote the expected cost associated with the admissible decision rule d in \mathcal{D} when the prior on H is p , i.e.,

$$J_p(d) \equiv \mathbb{E}_p[C(H, d(\mathbf{Y}))].$$

Similarly, let $J_p^*(\delta)$ denote the expected cost associated under the randomized decision rule δ in \mathcal{D}^* when the prior on H is p , i.e.,

$$J_p^*(\delta) \equiv \mathbb{E}_{\delta,p}[C(H, D)].$$

The Bayesian problems introduced in Chapters 2 and 3 now read

$$\mathcal{P}_{p,B} : \quad \text{Minimize } J_p(d) \text{ over } d \text{ in } \mathcal{D}$$

and

$$\mathcal{P}_{p,B}^* : \quad \text{Minimize } J_p^*(\delta) \text{ over } \delta \text{ in } \mathcal{D}^*.$$

The corresponding Bayesian costs will be denoted by

$$V(p) \equiv \inf_{d \in \mathcal{D}} J_p(d) \quad (4.2)$$

and

$$V^*(p) \equiv \inf_{\delta \in \mathcal{D}^*} J_p^*(\delta). \quad (4.3)$$

As shown in Chapter 2, under Condition **(A.1)**, for each p in $[0, 1]$ the problem $\mathcal{P}_{p,B}$ has a solution which we denote $d^*(p)$ to indicate its dependence on the prior p . Clearly, any such solution satisfies

$$J_p(d^*(p)) \leq J_p(d), \quad d \in \mathcal{D} \quad (4.4)$$

and the equality

$$V(p) = J_p(d^*(p)) \quad (4.5)$$

holds. Under the same condition, Theorem 3.4.1 further shows that

$$J_p^*(\delta_{d^*(p)}) \leq J_p(\delta), \quad \delta \in \mathcal{D}^*$$

so that

$$V^*(p) = V(p). \quad (4.6)$$

The following properties of the value function $V : [0, 1] \rightarrow \mathbb{R}$ will be useful in the forthcoming discussion. Conditions **(A.1)** and **(A.2)** are not needed for the results to hold.

Lemma 4.2.1 *Assume $\Gamma_h > 0$ for $h = 0, 1$. The value function $V : [0, 1] \rightarrow \mathbb{R}$ is concave and continuous on the closed interval $[0, 1]$ with boundary values $V(0) = C(0, 0)$ and $V(1) = C(1, 1)$. Moreover, its right-derivative (resp. left-derivative) exists and is finite on $[0, 1)$ (resp. $(0, 1]$)*

The proof can be omitted in a first reading, and can be found in Section 4.10. For easy reference, recall that for each p in $[0, 1]$ the expressions

$$\begin{aligned} J_p(d) &= pC(1, 1) + (1 - p)C(0, 0) \\ &\quad + \Gamma_0(1 - p) \cdot P_F(d) + \Gamma_1 p \cdot P_M(d), \quad d \in \mathcal{D} \end{aligned} \quad (4.7)$$

and

$$\begin{aligned} J_p^*(\delta) &= pC(1, 1) + (1 - p)C(0, 0) \\ &\quad + \Gamma_0(1 - p) \cdot P_F(\delta) + \Gamma_1 p \cdot P_M(\delta), \quad \delta \in \mathcal{D}^* \end{aligned} \quad (4.8)$$

hold. The relationships were given in Lemma 2.1.1 and Lemma 3.4.1, respectively.

4.3 The minimax formulation

Since the exact value of the prior p is not available, the Bayesian criterion has to be modified. Two different approaches are possible; each in its own way seeks to *compensate* for the uncertainty in the modeling assumptions.

Minimax One possible approach is to introduce a *worst-case* cost associated with the original cost, and then use it as the new criterion to be minimized. With this in mind, define

$$J_{\text{Max}}(d) \equiv \sup_{p \in [0,1]} J_p(d), \quad d \in \mathcal{D}. \quad (4.9)$$

We are then lead to consider the minimization problem

$$\mathcal{P}_{\text{Max}} : \quad \text{Minimize } J_{\text{Max}}(d) \text{ over } d \text{ in } \mathcal{D}.$$

Solving \mathcal{P}_{Max} amounts to finding an admissible test d_m^* in \mathcal{D} such that

$$J_{\text{Max}}(d_m^*) \leq J_{\text{Max}}(d), \quad d \in \mathcal{D}. \quad (4.10)$$

When it exists, the test d_m^* is known as a *minimax* test.

A priori there is no guarantee that a test in \mathcal{D} exists which satisfies (4.10) (even under Condition **(A.1)**) – It is not clear that a cost $\tilde{C} : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$ (likely related to the original cost $C : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$) and \tilde{p} in $[0, 1]$ can be found such that

$$J_{\text{Max}}(d) = \mathbb{E}_{\tilde{p}} \left[\tilde{C}(H, d(\mathbf{Y})) \right], \quad d \in \mathcal{D}.$$

If that were indeed the case, then Theorem 2.2.1 would guarantee the existence of a minimizer.

For technical reasons to become shortly apparent we also introduce the worst-case cost under randomized strategies, namely

$$J_{\text{Max}}^*(\delta) \equiv \sup_{p \in [0,1]} J_p^*(\delta), \quad \delta \in \mathcal{D}^*. \quad (4.11)$$

The minimization problem of interest here is now defined as

$$\mathcal{P}_{\text{Max}}^* : \quad \text{Minimize } J_{\text{Max}}^*(\delta) \text{ over } \delta \text{ in } \mathcal{D}^*.$$

Solving $\mathcal{P}_{\text{Max}}^*$ amounts to finding a randomized strategy δ_m^* in \mathcal{D}^* such that

$$J_{\text{Max}}^*(\delta_m^*) \leq J_{\text{Max}}^*(\delta), \quad \delta \in \mathcal{D}^*. \quad (4.12)$$

Again a priori there is no guarantee that there exists a test in \mathcal{D}^* satisfying (4.12) (even under Condition **(A.1)**). When it exists, the test δ_m^* is also known as a *minimax* test.

It is natural to wonder whether the tests d_m^* and δ_m^* exist, (possibly under additional conditions), whether they are different, and if not, whether $J_{\text{Max}}^*(\delta_m^*) = J_{\text{Max}}^*(d_m^*)$.

Maxmin Another reasonable way to proceed consists in using the Bayesian test for that value of p which yields the *largest* Bayesian cost (4.2): With the notation introduced earlier, let p_m in $[0, 1]$ such that

$$V(p_m) = \max_{p \in [0,1]} V(p), \quad (4.13)$$

and use the Bayesian rule $d^*(p_m)$ – The existence of p_m is guaranteed by the fact that the mapping $V : [0, 1] \rightarrow \mathbb{R}$ is continuous on the closed bounded interval $[0, 1]$ by Lemma 4.2.1, hence achieves its maximum value on $[0, 1]$.

The value p_m satisfying (4.13) is known as the *least favorable* prior. Although the terminology is not standard, we shall refer to $d^*(p_m)$ as a *maximin* test.

4.4 Preliminary facts

In view of the two competing approaches outlined in Section 4.3, several questions arise: (i) How does one characterize the minimax strategy d_m^* and develop ways find it; (ii) How does one characterize the least-favorable prior p_m and develop ways find it; (iii) Is there a simple relationship between the solutions proposed by two approaches, and in particular, whether is $d^*(p_m)$ is a candidate for d_m^* .

To frame the discussion of these issues we start with a couple of preliminary remarks.

The minimax inequalities As a first step towards understanding how the two approaches may be related to each other, consider the following arguments: From the definitions it always holds that

$$V(p) \leq J_p(d) \leq J_{\text{Max}}(d), \quad \begin{array}{l} p \in [0, 1] \\ d \in \mathcal{D}. \end{array} \quad (4.14)$$

It is now immediate that

$$V(p) \leq \inf_{d \in \mathcal{D}} J_{\text{Max}}(d), \quad p \in [0, 1]$$

since $V(p)$ does not depend on d , whence

$$\sup_{p \in [0, 1]} V(p) \leq \inf_{d \in \mathcal{D}} J_{\text{Max}}(d).$$

This last inequality can be rewritten as the *minimax inequality*

$$\sup_{p \in [0, 1]} \left(\inf_{d \in \mathcal{D}} J_p(d) \right) \leq \inf_{d \in \mathcal{D}} \left(\sup_{p \in [0, 1]} J_p(d) \right) \quad (4.15)$$

(in pure policies)

If we were to consider randomized strategies, it is also the case that

$$V^*(p) \leq J_p^*(\delta) \leq J_{\text{Max}}^*(\delta), \quad \begin{array}{l} p \in [0, 1] \\ \delta \in \mathcal{D}^* \end{array} \quad (4.16)$$

and arguments similar to the ones leading to (4.15) yield the minimax inequality

$$\sup_{p \in [0, 1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right) \leq \inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0, 1]} J_p^*(\delta) \right) \quad (4.17)$$

in randomized strategies.

Toward minimax equalities As we contrast the inequalities (4.15) and (4.17), it is natural to wonder whether these inequalities ever hold as *equalities*, namely

$$\sup_{p \in [0, 1]} \left(\inf_{d \in \mathcal{D}} J_p(d) \right) = \inf_{d \in \mathcal{D}} \left(\sup_{p \in [0, 1]} J_p(d) \right) \quad (4.18)$$

and

$$\sup_{p \in [0,1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right) = \inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0,1]} J_p^*(\delta) \right). \quad (4.19)$$

When this occurs we shall then say that the minimax property holds in pure and randomized policies, respectively.

It is worth pointing out that the equalities

$$\inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0,1]} J_p^*(\delta) \right) \leq \inf_{d \in \mathcal{D}} \left(\sup_{p \in [0,1]} J_p(d) \right) \quad (4.20)$$

and

$$\sup_{p \in [0,1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right) = \sup_{p \in [0,1]} \left(\inf_{d \in \mathcal{D}} J_p(d) \right) \quad (4.21)$$

always hold; the latter is a rewrite of (4.6) and is a simple consequence of Theorem 3.4.1. As we combine these observations with (4.17) we conclude that

$$\sup_{p \in [0,1]} \left(\inf_{d \in \mathcal{D}} J_p(d) \right) \leq \inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0,1]} J_p^*(\delta) \right) \leq \inf_{d \in \mathcal{D}} \left(\sup_{p \in [0,1]} J_p(d) \right). \quad (4.22)$$

Thus, if (4.18) happens to be true, then (4.19) necessarily holds – Put differently, the minimax property in pure policies is more difficult to achieve than the minimax property in randomized strategies. This disparity will become apparent in the discussion of the Minimax Theorem given in Section 4.5, opening the possibility that we may have to resort to randomized tests (at least in some situations) in order to achieve the minimax equality.

The structure of the worst-case costs (4.9) and (4.11) A little more can be said concerning the auxiliary costs (4.9) and (4.11): For each test d in \mathcal{D} , we note from (4.7) and (4.8) that

$$\sup_{p \in [0,1]} J_p(d) = \max_{p \in [0,1]} J_p(d) = \max\{J_0(d), J_1(d)\} \quad (4.23)$$

with the supremum achieved at either $p = 0$ or $p = 1$. Also, $J_0(d)$ and $J_1(d)$ can be given probabilistic interpretations as the conditional interpretations

$$J_0(d) = \mathbb{E}_p [C(H, d(\mathbf{Y}) | H = 0)] \quad (4.24)$$

and

$$J_1(d) = \mathbb{E}_p [C(H, d(\mathbf{Y})|H = 1)] \quad (4.25)$$

with p arbitrary in $[0, 1]$. Similarly, for each randomized strategy δ in \mathcal{D}^* , we have

$$\sup_{p \in [0,1]} J_p^*(\delta) = \max_{p \in [0,1]} J_p^*(\delta) = \max\{J_0^*(\delta), J_1^*(\delta)\} \quad (4.26)$$

with the supremum achieved at either $p = 0$ or $p = 1$ with probabilistic interpretations

$$J_0^*(\delta) = \mathbb{E}_{\delta,p} [C(H, D)|H = 0] \quad (4.27)$$

and

$$J_1^*(\delta) = \mathbb{E}_{\delta,p} [C(H, D)|H = 1] \quad (4.28)$$

with p arbitrary in $[0, 1]$.

4.5 The minimax equality

The main result concerning the minimax formulation for the binary hypothesis testing problem is summarized in the following special case of the Minimax Theorem from Statistical Decision Theory; see [?, Thm. 1, p. 82] for a discussion in a more general setting.

Theorem 4.5.1 *Assume $\Gamma_h > 0$ for all $h = 0, 1$. Under Condition (A.1), the minimax equality*

$$\sup_{p \in [0,1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right) = \inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0,1]} J_p^*(\delta) \right) \quad (4.29)$$

holds in randomized strategies.

In Section 4.6 we present an analysis of the minimax equality which exploits the specific structure of the binary hypothesis problem as reflected through the properties of the value function: As pointed out earlier, there always exists p_m in $[0, 1]$ such that (4.13) holds. From the concavity of the value function it follows that the set of maximizers

$$I_m \equiv \left\{ p_m \in [0, 1] : V(p_m) = \max_{p \in [0,1]} V(p) \right\}$$

is a closed interval in $[0, 1]$. The set I_m will often be reduced to a singleton, in which case the value function admits a unique (isolated) maximizer. Four situations can occur depending on the location of I_m and on the smoothness of $p \rightarrow V(p)$ at the maximum. In each case we establish a minimax equality and identify the minimax strategy. Throughout we still use $d^*(p_m)$ to denote the Bayesian test for the selected value p_m in I_m , so that

$$V(p_m) = J_{p_m}(d^*(p_m)) = \min_{d \in \mathcal{D}} J_{p_m}(d). \quad (4.30)$$

From the discussion of Section 4.4 we see that (4.29) will hold if we can establish the reverse inequality to (4.17), namely

$$\inf_{\delta \in \mathcal{D}^*} \left(\sup_{p \in [0,1]} J_p^*(\delta) \right) \leq \sup_{p \in [0,1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right). \quad (4.31)$$

Recall that (4.31) will automatically hold if we show the stronger inequality

$$\inf_{d \in \mathcal{D}} \left(\sup_{p \in [0,1]} J_p(d) \right) \leq \sup_{p \in [0,1]} \left(\inf_{d \in \mathcal{D}} J_p(d) \right). \quad (4.32)$$

In the first three cases we show in effect that

$$\inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right) \leq \max_{p \in [0,1]} \left(\min_{d \in \mathcal{D}} J_p(d) \right). \quad (4.33)$$

4.6 A proof of Theorem 4.5.1

We start with the boundary cases $p_m = 0$ and $p_m = 1$.

Case 1: Assume $p_m = 0$ – Thus, $\max_{p \in [0,1]} V(p) = V(0) = J_0(d^*(0))$. By concavity we have $\frac{d^+}{dp} V(p) \Big|_{p=0} \leq 0$ with the mapping $V : [0, 1] \rightarrow \mathbb{R}$ being decreasing. But the straight line $p \rightarrow J_p(d^*(0))$ is tangent to the value function $V : [0, 1] \rightarrow \mathbb{R}$ at $p = 0$, whence

$$\frac{d^+}{dp} V(p) \Big|_{p=0} = \frac{d}{dp} J_p(d^*(0)) \Big|_{p=0} \leq 0.$$

The mapping $p \rightarrow J_p(d^*(0))$ being affine, its derivative is therefore constant with

$$\frac{d}{dp} J_p(d^*(0)) = \frac{d}{dp} J_p(d^*(0)) \Big|_{p=0} \leq 0, \quad p \in [0, 1]$$

and the mapping $p \rightarrow J_p(d^*(0))$ is also decreasing on $[0, 1]$. This leads to

$$J_0(d^*(0)) = \max_{p \in [0,1]} J_p(d^*(0)).$$

With this in mind we get

$$\begin{aligned} \max_{p \in [0,1]} \left(\min_{d \in \mathcal{D}} J_p(d) \right) &= \max_{p \in [0,1]} V(p) \\ &= V(0) \\ &= J_0(d^*(0)) \\ &= \max_{p \in [0,1]} J_p(d^*(0)). \end{aligned} \quad (4.34)$$

The desired inequality (4.32) (hence (4.31)) is now immediate from (4.34) as we note that

$$\max_{p \in [0,1]} J_p(d^*(0)) \geq \inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right).$$

But the minimax equality being now established in pure strategies, we conclude from the discussion that

$$\max_{p \in [0,1]} J_p(d^*(0)) = \inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right).$$

This shows that d_m^* can be taken to be $d^*(0)$. ■

Case 2: Assume $p_m = 1$ – The proof is as in Case 1 *mutatis mutandis*, and is left as an exercise. Again, the minimax equality holds in pure strategies and d_m^* can be taken to be $d^*(1)$. ■

We now turn to cases when p_m is selected in $(0, 1)$.

Case 3: Assume that p_m in an element of $(0, 1)$ and $p \rightarrow V(p)$ is differentiable at $p = p_m$ – It is plain that $\left. \frac{d}{dp} V(p) \right|_{p=p_m} = 0$ since p_m is an interior point by assumption. By concavity the mapping $p \rightarrow J_p(d^*(p_m))$ is tangent to the value function $V : [0, 1] \rightarrow \mathbb{R}$ at $p = p_m$, whence

$$\left. \frac{d}{dp} V(p) \right|_{p=p_m} = \left. \frac{d}{dp} J_p(d^*(p_m)) \right|_{p=p_m} = 0.$$

The mapping $p \rightarrow J_p(d^*(p_m))$ being affine, its derivative is constant and given by

$$\frac{d}{dp} J_p(d^*(p_m)) = \frac{d}{dp} J_p(d^*(p_m)) \Big|_{p=0} = 0, \quad p \in [0, 1].$$

Therefore, the mapping $p \rightarrow J_p(d^*(p_m))$ is constant on $[0, 1]$, and the equality $J_0(d^*(p_m)) = J_1(d^*(p_m))$ holds. It follows from the first equality in (4.30) that

$$V(p_m) = J_p(d^*(p_m)) = \max_{p \in [0,1]} J_p(d^*(p_m)), \quad p \in [0, 1]. \quad (4.35)$$

On the other hand, it is plain that

$$\begin{aligned} \inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right) &\leq \max_{p \in [0,1]} J_p(d^*(p_m)) \\ &= J_{p_m}(d^*(p_m)) \\ &= \min_{d \in \mathcal{D}} J_{p_m}(d) \\ &\leq \inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right) \end{aligned} \quad (4.36)$$

as we use the second equality in (4.35) with $p = p_m$, and then apply the second equality in (4.30). The inequality (4.32) (hence (4.31)) is now a straightforward consequence of (4.36).

Leveraging the fact that the minimax equality is now known to hold in pure strategies, we conclude from the discussion that

$$\max_{p \in [0,1]} J_p(d^*(p_m)) = \inf_{d \in \mathcal{D}} \left(\max_{p \in [0,1]} J_p(d) \right),$$

and d_m^* can therefore be taken to be $d^*(p_m)$. ■

Case 4: Assume that $I_m = \{p_m\} \subseteq (0, 1)$ but $p \rightarrow V(p)$ is not differentiable at $p = p_m$ – Under such assumptions we must have

$$a_+ \equiv \frac{d^+}{dp} V(p) \Big|_{p=p_m} < \frac{d^-}{dp} V(p) \Big|_{p=p_m} \equiv a_-$$

by concavity with either $a_+ < 0 \leq a_-$ or $a_+ \leq 0 < a_-$. We continue the discussion under the assumption $a_+ < 0 \leq a_-$; the case $a_+ \leq 0 < a_-$ proceeds along similar lines, and is therefore omitted.

Recall that $p \rightarrow V(p)$ is defined as the envelope of a family of affine functions. Thus, under the non-differentiability assumption at $p = p_m$, concavity guarantees that there exist two pure strategies, say $d_-, d_+ : \mathbb{R}^k \rightarrow \{0, 1\}$, such that $V(p_m) = J_{p_m}(d_-)$ and $V(p_m) = J_{p_m}(d_+)$ (because p_m is a maximum) while the straight lines $p \rightarrow J_p(d_-)$ and $p \rightarrow J_p(d_+)$ are both tangent to the value function at $p = p_m$ – These two strategies are distinct. Hence, as discussed in earlier cases, the function $p \rightarrow J_p(d_-)$ (resp. $p \rightarrow J_p(d_+)$) is an affine function with constant derivative $a_- \geq 0$ (resp. $a_+ < 0$), hence non-decreasing (resp. decreasing). It follows that $J_0(d_-) \leq J_1(d_-)$ and $J_1(d_+) < J_0(d_+)$.

Next we introduce randomized policies $\{\delta_a, a \in [0, 1]\}$ obtained by randomizing two pure strategies d_- and d_+ . Thus, with each a in $[0, 1]$ consider the randomized policy $\delta_a : \mathbb{R}^k \rightarrow [0, 1]$ given by

$$\delta_a = ad_+ + (1 - a)d_-.$$

The relation (3.20) discussed in Section 3.5 applies, yielding

$$J_p^*(\delta_a) = aJ_p(d_+) + (1 - a)J_p(d_-), \quad p \in [0, 1].$$

By construction we also note that

$$V(p_m) = J_{p_m}^*(\delta_a), \quad a \in [0, 1]. \quad (4.37)$$

If a suitable of a , we were to have $p \rightarrow J_p^*(\delta_a)$ constant over $[0, 1]$, then the test δ_a would a performance insensitive to the value of p . This requirement (on a) is equivalent to the equality $J_0^*(\delta_a) = J_1^*(\delta_a)$, i.e.,

$$aJ_0(d_+) + (1 - a)J_0(d_-) = aJ_1(d_+) + (1 - a)J_1(d_-).$$

Thus,

$$a((J_0(d_+) - J_1(d_+)) + (J_1(d_-) - J_0(d_-))) = J_1(d_-) - J_0(d_-)$$

and solving for a we get

$$a^* = \frac{J_1(d_-) - J_0(d_-)}{(J_0(d_+) - J_1(d_+)) + (J_1(d_-) - J_0(d_-))}.$$

It is a simple matter to check that a^* lies in $[0, 1)$ since $J_1(d_-) - J_0(d_-) \geq 0$ and $J_0(d_+) - J_1(d_+) > 0$ as discussed earlier.

It is now plain that

$$V(p_m) = J_p^*(\delta_{a^*}) = \max_{p \in [0,1]} J_p^*(\delta_{a^*}), \quad p \in [0, 1]. \quad (4.38)$$

Therefore, as in the discussion for Case 3, we have

$$\begin{aligned} \inf_{\delta \in \mathcal{D}^*} \left(\max_{p \in [0,1]} J_p^*(\delta) \right) &\leq \max_{p \in [0,1]} J_p^*(\delta_{a^*}) \\ &= V(p_m) \\ &= \inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \\ &\leq \max_{p \in [0,1]} \left(\inf_{\delta \in \mathcal{D}^*} J_p^*(\delta) \right) \end{aligned} \quad (4.39)$$

and the desired inequality (4.31) is established.

The minimax equality now holds in mixed strategies, whence

$$\max_{p \in [0,1]} J_p^*(\delta_{a^*}) = \inf_{\delta \in \mathcal{D}^*} \left(\max_{p \in [0,1]} J_p^*(\delta) \right)$$

by virtue of (4.39). The minimax strategy is a randomized strategy δ_m^* which is identified as $\delta^*(p_m)$. Note that $\delta^*(p_m)$ is also a (randomized) Bayesian policy for the least favorable prior. ■

We summarize these findings in the following corollary to Theorem 4.5.1.

Corollary 4.6.1 *Assume $\Gamma_h > 0$ for all $h = 0, 1$. Under Condition (A.1), the minimax equality (4.29) holds in randomized strategies. Moreover, the minimax strategy always exists and can be interpreted as a (possibly randomized) Bayesian test under the least favorable prior p_m .*

4.7 The minimax equation

The discussion of Section 4.5 shows that finding minimax tests passes through the evaluation of the value function $p \rightarrow V(p)$ and its maximizing set I_m . As

simple examples already suggest in later sections, this evaluation may not always be straightforward to carry. Moreover, once an expression for the value function becomes available, finding its maximizers may turn out to be rather cumbersome. However, this two-step approach can be bypassed when I_m contains an interior point p_m at which the value function is differentiable, in which case the minimax test is given by the Bayesian test $d^*(p_m)$. Instead a simple characterization of p_m is achieved through the so-called *Minimax Equation*.

Lemma 4.7.1 *Assume $\Gamma_h > 0$ for all $h = 0, 1$, and that p_m is an element of $(0, 1)$ and $p \rightarrow V(p)$ is differentiable at $p = p_m$. Under Condition (A.1), p_m can be characterized through the Minimax Equation*

$$C(1, 1) - C(0, 0) = \Gamma_0 \cdot P_F(d^*(p_m)) - \Gamma_1 \cdot P_M(d^*(p_m)). \quad (4.40)$$

For the probability of error criterion, the Minimax Equation takes the simpler form

$$P_F(d^*(p_m)) = P_M(d^*(p_m)). \quad (4.41)$$

Proof. Fix p in $[0, 1]$. Upon specializing (4.7) to the test $d^*(p)$, we get

$$\begin{aligned} J_\alpha(d^*(p)) &= \alpha C(1, 1) + (1 - \alpha)C(0, 0) \\ &\quad + \Gamma_0(1 - \alpha) \cdot P_F(d^*(p)) + \Gamma_1 \alpha \cdot P_M(d^*(p)) \end{aligned} \quad (4.42)$$

with α in $[0, 1]$ and the mapping $\alpha \rightarrow J_\alpha(d^*(p))$ is therefore affine in the variable α on the interval $[0, 1]$. Therefore, the graph of the mapping $\alpha \rightarrow J_\alpha(d^*(p))$ is a straight line; its slope is given by

$$\frac{d}{d\alpha} J_\alpha(d^*(p)) = C(1, 1) - C(0, 0) + \Gamma_1 \cdot P_M(d^*(p)) - \Gamma_0 \cdot P_F(d^*(p)). \quad (4.43)$$

By its definition, the Bayesian cost satisfies

$$V(\alpha) \leq J_\alpha(d), \quad \begin{array}{l} d \in \mathcal{D} \\ \alpha \in [0, 1] \end{array}$$

with strict inequality for most tests. With $d = d^*(p)$ this inequality becomes an equality when $\alpha = p$, namely

$$V(p) = J_p(d^*(p))$$

while

$$V(\alpha) \leq J_\alpha(d^*(p)), \quad \alpha \in [0, 1].$$

With p in $(0, 1)$, if the concave mapping $\alpha \rightarrow V(\alpha)$ is *differentiable* at $\alpha = p$, then the straight line $\alpha \rightarrow J_\alpha(d^*(p))$ will be a tangent to the mapping $\alpha \rightarrow V(\alpha)$ at $\alpha = p$ – This is a consequence of the concavity established in Lemma 4.2.1. Thus,

$$\left. \frac{d}{d\alpha} V(\alpha) \right|_{\alpha=p} = \left. \frac{d}{d\alpha} J_\alpha(d^*(p)) \right|_{\alpha=p}. \quad (4.44)$$

In particular, if p_m is an element of $(0, 1)$ and the mapping $\alpha \rightarrow V(\alpha)$ is differentiable at $\alpha = p_m$, then

$$\left. \frac{d}{d\alpha} V(\alpha) \right|_{\alpha=p_m} = \left. \frac{d}{d\alpha} J_\alpha(d^*(p_m)) \right|_{\alpha=p_m}. \quad (4.45)$$

But the interior point p_m being a maximum for the function $\alpha \rightarrow V(\alpha)$, we must have

$$\left. \frac{d}{d\alpha} V(\alpha) \right|_{\alpha=p_m} = 0,$$

whence

$$\left. \frac{d}{d\alpha} J_\alpha(d^*(p_m)) \right|_{\alpha=p_m} = 0.$$

The equation (4.40) now follows from (4.43). ■

Obviously this analysis does not cover the cases when (i) $p_m = 0$, (ii) $p_m = 1$ and (iii) p_m is an element of $(0, 1)$ but the mapping $\alpha \rightarrow V(\alpha)$ is not differentiable at $\alpha = p_m$.

4.8 The Gaussian Case

The setting is that of Section 2.4 to which we refer the reader for the notation. As shown there, for every $\eta > 0$ we have

$$P_F(Lrt_\eta) = 1 - \Phi \left(\frac{\log \eta + \frac{1}{2}d^2}{d} \right)$$

and

$$P_M(Lrt_\eta) = \Phi \left(\frac{\log \eta - \frac{1}{2}d^2}{d} \right).$$

For each p in $(0, 1]$, with

$$\eta(p) = \frac{1-p}{p} \cdot \frac{\Gamma_0}{\Gamma_1},$$

we have $d^*(p) = Lrt_{\eta(p)}$ and the expression (4.7) yields

$$\begin{aligned} V(p) &= J_p(d^*(p)) \\ &= pC(1, 1) + (1-p)C(0, 0) \\ &\quad + \Gamma_0(1-p) \cdot \left(1 - \Phi\left(\frac{\log \eta(p) + \frac{1}{2}d^2}{d}\right)\right) \\ &\quad + \Gamma_1 p \cdot \Phi\left(\frac{\log \eta(p) - \frac{1}{2}d^2}{d}\right). \end{aligned} \quad (4.46)$$

The boundary cases $p = 0$ is easily recovered upon formally substituting this value in the expression (4.46). The Minimax Equation (4.40) takes the form

$$\begin{aligned} C(1, 1) - C(0, 0) \\ = \Gamma_1 \Phi\left(\frac{\log \eta(p_m) - \frac{1}{2}d^2}{d}\right) - \Gamma_0 \left(1 - \Phi\left(\frac{\log \eta(p_m) + \frac{1}{2}d^2}{d}\right)\right). \end{aligned} \quad (4.47)$$

Probability of error – Simplifications occur since $C(0, 0) = C(1, 1) = 0$ and $\Gamma_0 = \Gamma_1 = 1$: The expression (4.46) becomes

$$V(p) = (1-p) \cdot \left(1 - \Phi\left(\frac{\frac{1-p}{p} + \frac{1}{2}d^2}{d}\right)\right) + p \cdot \Phi\left(\frac{\log \frac{1-p}{p} - \frac{1}{2}d^2}{d}\right),$$

and the Minimax Equation (4.47) reduces to

$$\Phi\left(\frac{\log \eta(p_m) - \frac{1}{2}d^2}{d}\right) + \Phi\left(\frac{\log \eta(p_m) + \frac{1}{2}d^2}{d}\right) = 1.$$

It is easy to see that this requires $\log \eta(p_m) = 0$ so that $p_m = \frac{1}{2}$ (indeed in $(0, 1)$), an intuitively satisfying conclusion! Moreover, the minimax test is given by $d_m^* = d(\frac{1}{2})$.

4.9 The Bernoulli case

The setting is that of Section 2.5 to which we refer the reader for the notation. We discuss only the case $a_1 < a_0$, and leave the case $a_0 < a_1$ as an exercise for the interested reader.

Note that the condition $a_1 < a_0$ is equivalent to $1 < \frac{1-a_1}{1-a_0}$, so that the expressions (2.30) and (2.31) for the probabilities $P_F(d_\eta)$ and $P_M(d_\eta)$, respectively, are *piecewise* constant functions of η with different constant values on the intervals $(0, \frac{a_1}{a_0}]$, $(\frac{a_1}{a_0}, \frac{1-a_1}{1-a_0}]$ and $(\frac{1-a_1}{1-a_0}, \infty)$: Direct inspection of the expression (2.30) yields

$$P_F(d_\eta) = \begin{cases} 1 & \text{if } 0 < \eta \leq \frac{a_1}{a_0} \\ 1 - a_0 & \text{if } \frac{a_1}{a_0} < \eta \leq \frac{1-a_1}{1-a_0} \\ 0 & \text{if } \frac{1-a_1}{1-a_0} < \eta. \end{cases} \quad (4.48)$$

Similarly, using (2.31) we find

$$P_M(d_\eta) = \begin{cases} 0 & \text{if } 0 < \eta \leq \frac{a_1}{a_0} \\ a_1 & \text{if } \frac{a_1}{a_0} < \eta \leq \frac{1-a_1}{1-a_0} \\ 1 & \text{if } \frac{1-a_1}{1-a_0} < \eta. \end{cases} \quad (4.49)$$

Thus, for each p in $[0, 1]$, we see from (4.7) that the cost $J_p(d_\eta)$ takes a different value on each of the intervals $(0, \frac{a_1}{a_0}]$, $(\frac{a_1}{a_0}, \frac{1-a_1}{1-a_0}]$ and $(\frac{1-a_1}{1-a_0}, \infty)$: Specifically, we have:

On $(0, \frac{a_1}{a_0}]$,

$$\begin{aligned} J_p(d_\eta) &= pC(1, 1) + (1-p)C(0, 0) + \Gamma_0(1-p) \\ &= pC(1, 1) + (1-p)C(0, 1). \end{aligned} \quad (4.50)$$

On $(\frac{a_1}{a_0}, \frac{1-a_1}{1-a_0}]$,

$$\begin{aligned} J_p(d_\eta) &= pC(1, 1) + (1-p)C(0, 0) + \Gamma_0(1-p) \cdot (1-a_0) + \Gamma_1 p \cdot a_1 \\ &= pC(1, 1) + (1-p)C(0, 1) + \Gamma_1 p \cdot a_1 - \Gamma_0(1-p) \cdot a_0 \\ &= p(C(1, 1) + \Gamma_1 a_1) + (1-p)(C(0, 1) - \Gamma_0 a_0). \end{aligned} \quad (4.51)$$

On $(\frac{1-a_1}{1-a_0}, \infty)$,

$$\begin{aligned} J_p(d_\eta) &= pC(1, 1) + (1-p)C(0, 0) + \Gamma_1 p \\ &= pC(1, 0) + (1-p)C(0, 0). \end{aligned} \quad (4.52)$$

Recall that

$$V(p) = J_p(d_{\eta(p)}) \text{ with } \eta(p) = \frac{\Gamma_0(1-p)}{\Gamma_1 p}, \quad 0 < p \leq 1.$$

As the mapping $p : (0, 1] \rightarrow \mathbb{R}_+ : p \rightarrow \eta(p)$ is strictly decreasing, each of the equations

$$\eta(p) = \frac{1-a_1}{1-a_0}, \quad 0 < p \leq 1$$

and

$$\eta(p) = \frac{a_1}{a_0}, \quad 0 < p \leq 1$$

has a unique solution in $(0, 1)$. These solutions, denoted p_- and p_+ , respectively, are given by

$$p_- = \frac{\Gamma_0(1-a_0)}{\Gamma_1(1-a_1) + \Gamma_0(1-a_0)}$$

and

$$p_+ = \frac{\Gamma_0 a_0}{\Gamma_1 a_1 + \Gamma_0 a_0}.$$

As expected $p_- < \frac{1}{2} < p_+$.

Earlier expressions can now be used, and yield

$$V(p) = \begin{cases} pC(1, 0) + (1-p)C(0, 0) & \text{if } p \in (0, p_-] \\ p(C(1, 1) + \Gamma_1 a_1) + (1-p)(C(0, 1) - \Gamma_0 a_0) & \text{if } p \in (p_-, p_+] \\ pC(1, 1) + (1-p)C(0, 1) & \text{if } p \in (p_+, 1). \end{cases}$$

It is plain that the function $V : [0, 1] \rightarrow \mathbb{R}$ is piecewise linear with three distinct segments, namely $(0, p_-]$, $(p_-, p_+]$ and $(p_+, 1]$. There are two kinks at $p = p_-$ and $p = p_+$, respectively. That the function is concave can be seen by computing the left and right-derivatives at these points. The function $V : [0, 1] \rightarrow \mathbb{R}$ is differentiable everywhere except at these kinks. However the maximum occurs at one of these points so that $p_m \in \{p_-, p_+\}$.

Probability of error – In that case we find that

$$V(p) = \begin{cases} p & \text{if } p \in (0, p_-] \\ pa_1 + (1-p)(1-a_0) & \text{if } p \in (p_-, p_+] \\ 1-p & \text{if } p \in (p_+, 1) \end{cases} \quad (4.53)$$

with

$$p_- = \frac{1-a_0}{(1-a_1)+(1-a_0)}$$

and

$$p_+ = \frac{a_0}{a_1+a_0}.$$

It is a simple matter to check that $V(p_{\pm-}) = V(p_{\pm+})$, establishing continuity at the kink points. As we compare $V(p_-)$ and $V(p_+)$, we readily conclude that $p_m = p_-$ (resp. $p_m = p_+$) iff $1-p_+ < p_-$ (resp. $p_- < 1-p_+$) iff $a_0 + a_1 < 1$ (resp. $1 < a_0 + a_1$). The minimax cost is then given by

$$V_m = \begin{cases} p_- = \frac{1-a_0}{(1-a_1)+(1-a_0)} & \text{if } a_0 + a_1 < 1 \\ 1-p_+ = \frac{a_1}{a_1+a_0} & \text{if } 1 < a_0 + a_1 \end{cases}$$

Minimax strategy is necessarily randomized and is given by

$$\delta_a = ad_+ + (1-a)d_-$$

with the pure tests $d_-, d_+ : \mathbb{R} \rightarrow \{0, 1\}$ given by

4.10 A proof of Lemma 4.2.1

The proof proceeds in several stages. We start with the fact that

$$V(p) = \inf_{d \in \mathcal{D}} J_p(d), \quad p \in [0, 1].$$

Values at the boundary points – Consider a test d in \mathcal{D} . With $p = 0$ and $p = 1$ in (4.7) we get

$$J_0(d) = C(0, 0) + \Gamma_0 P_F(d)$$

and

$$J_1(d) = C(1, 1) + \Gamma_1 P_M(d).$$

Using the conditions $\Gamma_0 > 0$ and $\Gamma_1 > 0$, we conclude that

$$V(0) = \inf_{d \in \mathcal{D}} J_0(d) = C(0, 0) + \Gamma_0 \cdot \inf_{d \in \mathcal{D}} P_F(d)$$

and

$$V(1) = \inf_{d \in \mathcal{D}} J_1(d) = C(1, 1) + \Gamma_1 \cdot \inf_{d \in \mathcal{D}} P_M(d).$$

However, $P_F(d_F) = 0$ for the test $d_F : \mathbb{R}^k \rightarrow \{0, 1\}$ which always selects the null hypothesis ($H = 0$) while $P_M(d_M) = 0$ for the test $d_M : \mathbb{R}^k \rightarrow \{0, 1\}$ which always selects the alternative ($H = 1$). It follows that $\inf_{d \in \mathcal{D}} P_F(d) = 0$ and $\inf_{d \in \mathcal{D}} P_M(d) = 0$, whence $V(0) = C(0, 0)$ and $V(1) = C(1, 1)$. ■

Concavity on $[0, 1]$ and continuity on $(0, 1)$ – Once the test d is selected, the probabilities $P_F(d)$ and $P_M(d)$ appearing in (4.7) do *not* depend on p , and are determined only through F_0 and F_1 . Thus, the mapping $p \rightarrow J_p(d)$ is affine, hence concave in p . As a result, the mapping $V : [0, 1] \rightarrow \mathbb{R}$ is concave on the closed interval $[0, 1]$, being the infimum of the family $\{J_p(d), d \in \mathcal{D}\}$ of concave functions. Because a concave function defined on an open interval is necessarily continuous on that open interval, the mapping $V : [0, 1] \rightarrow \mathbb{R}$ is continuous on $(0, 1)$ by virtue of Fact 9.4.2. ■

Continuity at the boundary points – We now turn to showing that the mapping $V : [0, 1] \rightarrow \mathbb{R}$ is also continuous at the boundary points $p = 0$ and $p = 1$. We discuss only the case $p = 0$; the case $p = 1$ can be handled *mutatis mutandis* and is left to the interested reader as an exercise.

For notational convenience here and below we write

$$\Delta(p) \equiv \inf_{d \in \mathcal{D}} (\Gamma_0(1 - p) \cdot P_F(d) + \Gamma_1 p \cdot P_M(d)), \quad p \in (0, 1].$$

Recall that $V(0) = C(0, 0)$ by the first part of the proof. Thus, for each p in $(0, 1]$ we get from the definition of $V(p)$ that

$$V(p) - V(0) = p(C(1, 1) - C(0, 0)) + \Delta(p) \quad (4.54)$$

by virtue of (4.7). The continuity of the mapping $V : [0, 1] \rightarrow \mathbb{R}$ at $p = 0$ is therefore equivalent to

$$\lim_{p \rightarrow 0} \Delta(p) = 0. \quad (4.55)$$

For any fixed p in $(0, 1]$, the conditions $\Gamma_0 > 0$ and $\Gamma_1 > 0$ yield the inequalities

$$0 \leq \Delta(p) \leq \Gamma_1 p \quad (4.56)$$

since under the test d_F (introduced earlier in the proof) we have $P_F(d_F) = 0$ and $P_M(d_F) = 1$. The conclusion (4.55) is now immediate. ■

Differentiability – The existence *and* finiteness of the right-derivative and left-derivative on the open interval $(0, 1)$ are simple consequences of Fact 9.4.4. The same argument also shows that the right-derivative (resp. left-derivative) does exist at $p = 0$ (resp. $p = 1$); however it may not necessarily be finite.

Instead, we provide a direct argument to show the existence and finiteness of the right-derivative (resp. left-derivative) at $p = 0$ (resp. $p = 1$). We carry out the discussion only for $p = 0$ as the case $p = 1$ is similar: For each p in $(0, 1]$, we note that

$$\frac{V(p) - V(0)}{p} = C(1, 1) - C(0, 0) + \frac{\Delta(p)}{p} \quad (4.57)$$

with

$$\frac{\Delta(p)}{p} = \inf_{d \in \mathcal{D}} \left(\Gamma_0 \left(\frac{1}{p} - 1 \right) \cdot P_F(d) + \Gamma_1 \cdot P_M(d) \right).$$

This last expression shows that $p \rightarrow \frac{\Delta(p)}{p}$ is decreasing on $(0, 1]$, whence the limit $\lim_{p \downarrow 0} \frac{\Delta(p)}{p}$ always exists. This limit is finite by virtue of the bounds

$$0 \leq \frac{\Delta(p)}{p} \leq \Gamma_1, \quad p \in (0, 1]$$

which are inherited from the earlier bounds (4.56). This shows the existence of a finite right-derivative at $p = 0$. ■

4.11 Exercises

4.1.

Let I denote an interval of \mathbb{R} , not necessarily finite, closed or open, and let A be an arbitrary index set. For each α in A , let $f_\alpha : I \rightarrow \mathbb{R}$ be a concave function. With the function $g : I \rightarrow \mathbb{R}$ defined by

$$g(x) = \inf (f_\alpha(x) : \alpha \in A), \quad x \in I$$

show that the mapping $g : I \rightarrow \mathbb{R}$ is concave.

4.2.

With $h > 0$ show that the equation

$$\Phi(x - h) + \Phi(x + h) = 1, \quad x \in \mathbb{R}$$

has a unique solution $x = 0$.

4.3.

4.12 References

Chapter 5

The Neyman-Pearson formulation

In many situations, not only is the prior probability p not available but it is quite difficult to make meaningful cost assignments. This is typically the case in radar applications – After all, what is the real cost of failing to detect an incoming missile? While it is tempting to seek to minimize *both* the probabilities of miss and false alarm, these are (usually) conflicting objectives and a *constrained* optimization problem is considered instead. The Neyman-Pearson formulation of the binary hypothesis problem given next constitutes an approach to handle such situations.

5.1 A constrained optimization problem

Fix α in $(0, 1)$ (the limiting case $\alpha = 0$ being of little practical interest). Let \mathcal{D}_α denote the collection of admissible tests in \mathcal{D} of *size* at most α , namely

$$\mathcal{D}_\alpha = \{d \in \mathcal{D} : P_F(d) \leq \alpha\}.$$

The Neyman-Pearson formulation is based on solving the constrained optimization problem NP_α where

$$\text{NP}_\alpha : \quad \text{Maximize } P_D(d) \text{ over } d \text{ in } \mathcal{D}_\alpha.$$

Solving NP_α amounts to finding a test $d_{\text{NP}}(\alpha)$ in \mathcal{D}_α with the property that

$$P_D(d) \leq P_D(d_{\text{NP}}(\alpha)), \quad d \in \mathcal{D}_\alpha.$$

Such a test $d_{\text{NP}}(\alpha)$, when it exists, is called a *Neyman-Pearson* test of size α , or

alternatively, an α -level Neyman–Pearson decision rule. Such decision rules may not be unique. Following the accepted terminology, its *power* $\beta(\alpha)$ is given by

$$\beta(\alpha) \equiv P_D(d_{\text{NP}}(\alpha)) = \sup_{d \in \mathcal{D}_\alpha} P_D(d).$$

When reformulated as

$$\text{NP}_\alpha : \quad \text{Minimize } P_M(d) \text{ over } d \text{ in } \mathcal{D}_\alpha,$$

the constrained optimization problem NP_α can be solved by standard Lagrangian arguments which are outlined in the next section. Throughout we assume that Condition (A.1) holds.

5.2 The Lagrangian arguments

Fix α in $(0, 1)$. For each $\lambda \geq 0$ consider the Lagrangian functional $J_\lambda : \mathcal{D} \rightarrow \mathbb{R}$ given by

$$J_\lambda(d) = P_M(d) + \lambda(P_F(d) - \alpha), \quad d \in \mathcal{D}.$$

The *Lagrangian* problem LP_λ is now defined as the *unconstrained* minimization problem

$$\text{LP}_\lambda : \quad \text{Minimize } J_\lambda(d) \text{ over } d \text{ in } \mathcal{D}.$$

Solving LP_λ amounts to finding a test d_λ^* in \mathcal{D} such that

$$J_\lambda(d_\lambda^*) \leq J_\lambda(d), \quad d \in \mathcal{D}.$$

Solving the Lagrangian problem LP_λ Fix $\lambda > 0$. For any test d in \mathcal{D} , we note that

$$\begin{aligned} J_\lambda(d) &= \mathbb{P}[d(\mathbf{Y}) = 0 | H = 1] + \lambda(\mathbb{P}[d(\mathbf{Y}) = 1 | H = 0] - \alpha) \\ &= \mathbb{P}[d(\mathbf{Y}) = 0 | H = 1] + \lambda(1 - \mathbb{P}[d(\mathbf{Y}) = 0 | H = 0] - \alpha) \\ &= \lambda(1 - \alpha) + \mathbb{P}[d(\mathbf{Y}) = 0 | H = 1] - \lambda\mathbb{P}[d(\mathbf{Y}) = 0 | H = 0] \\ &= \lambda(1 - \alpha) + \int_{C(d)} h_\lambda(\mathbf{y}) dF(\mathbf{y}) \end{aligned} \tag{5.1}$$

with $h_\lambda : \mathbb{R}^k \rightarrow \mathbb{R}$ given by

$$h_\lambda(\mathbf{y}) = f_1(\mathbf{y}) - \lambda f_0(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k.$$

By the comparison arguments used in the proof of Theorem 2.2.1, the Lagrangian problem LP_λ is easily seen to be solved by the test $d_\lambda^* : \mathbb{R}^k \rightarrow \{0, 1\}$ given by

$$d_\lambda^*(\mathbf{y}) = 0 \quad \text{iff} \quad h_\lambda(\mathbf{y}) < 0, \quad (5.2)$$

or equivalently,

$$d_\lambda^*(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < \lambda f_0(\mathbf{y}). \quad (5.3)$$

Note that in the notation associated with the definition (??) we have $d_\lambda^* = d_\lambda$. ■

Meeting the constraint The next step consists in finding some value $\lambda(\alpha) > 0$ of the Lagrangian multiplier such that the test $d_{\lambda(\alpha)}$ meets the constraint, i.e.,

$$P_F(d_{\lambda(\alpha)}) = \alpha. \quad (5.4)$$

If such value $\lambda(\alpha)$ were to exist, then the optimality $d_{\lambda(\alpha)}$ implies

$$J_{\lambda(\alpha)}(d_{\lambda(\alpha)}) \leq J_{\lambda(\alpha)}(d), \quad d \in \mathcal{D},$$

or equivalently,

$$P_M(d_{\lambda(\alpha)}) \leq P_M(d) + \lambda(\alpha) (P_F(d) - \alpha), \quad d \in \mathcal{D}.$$

Consequently, for every test d in \mathcal{D}_α (and not merely in \mathcal{D}), it follows that

$$P_M(d_{\lambda(\alpha)}) \leq P_M(d)$$

since then $P_M(d) \leq \alpha$. The test $d_{\lambda(\alpha)}$ is a test in \mathcal{D}_α by virtue of (5.4), hence it solves NP_α – In other words, $d_{NP}(\alpha)$ can be taken to be $d_{\lambda(\alpha)}$. ■

A difficulty The Lagrangian argument hinges upon the possibility of finding a value $\lambda(\alpha)$ of the Lagrange multiplier such that $P_F(d_{\lambda(\alpha)}) = \alpha$. Unfortunately, this may not be always possible, unless additional assumptions are imposed. To see how this may indeed happen, note that

$$\begin{aligned} P_F(d_\lambda) &= \mathbb{P}[d_\lambda(\mathbf{Y}) = 1 | H = 0] \\ &= \mathbb{P}[f_1(\mathbf{Y}) \geq \lambda f_0(\mathbf{Y}) | H = 0], \quad \lambda > 0. \end{aligned} \quad (5.5)$$

The mapping $\mathbb{R}_+ \rightarrow [0, 1] : \lambda \rightarrow P_F(d_\lambda)$ is clearly *monotone non-increasing*. However, the constraint $P_F(d_\lambda) = \alpha$ may fail to hold for some α in $(0, 1]$ because the set of values $\{P_F(d_\lambda), \lambda \geq 0\}$ need not contain α . This will occur if the mapping $\lambda \rightarrow P_F(d_\lambda)$ is not continuous at some point, say $\lambda^* > 0$, with

$$\lim_{\lambda \uparrow \lambda^*} P_F(d_\lambda) < \alpha < \lim_{\lambda \downarrow \lambda^*} P_F(d_\lambda).$$

In Section ?? we illustrate such situations on simple examples that involve discrete rvs. Randomized policies are introduced to solve this difficulty. There are however situations where this can be avoided because each one of the problems NP_α (properly defined over randomized strategies) has a solution within the set of non-randomized policies \mathcal{D} .

5.3 The Neyman-Pearson Lemma

The discussion of Section 5.2 suggests the need to consider an extended version of the Neyman-Pearson formulation where randomized strategies are allowed.

Fix α in $(0, 1]$. Let \mathcal{D}_α^* denote the collection of all randomized tests in \mathcal{D}^* of size at most α , namely

$$\mathcal{D}_\alpha^* = \{\delta \in \mathcal{D}^* : P_F(\delta) \leq \alpha\}.$$

The constrained optimization problem NP^* is now replaced by the following constrained optimization problem NP_α^* where

$$\text{NP}_\alpha^* : \quad \text{Maximize } P_D(\delta) \text{ over } \delta \text{ in } \mathcal{D}_\alpha^*.$$

Solving NP_α^* amounts to finding a test $\delta_{\text{NP}}(\alpha)$ in \mathcal{D}_α^* with the property that

$$P_D(\delta) \leq P_D(\delta_{\text{NP}}(\alpha)), \quad \delta \in \mathcal{D}_\alpha^*.$$

Such a test $\delta_{\text{NP}}(\alpha)$, when it exists, is also called a *Neyman–Pearson* test of size α , or alternatively, an α -level Neyman–Pearson decision rule. It may not be unique.

The existence of the Neyman–Pearson test $\delta_{\text{NP}}(\alpha)$ of size α , its characterization and uniqueness are discussed below through three separate lemmas, known collectively as the Neyman-Pearson Lemma. Proofs are delayed until Section 5.4.

First a definition: With $\eta \geq 0$ and Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ (to be selected shortly), define the randomized test $\delta^* : \mathbb{R}^k \rightarrow [0, 1]$ in \mathcal{D}^* given by

$$\delta^*(\mathbf{y}) = \begin{cases} 1 & \text{if } \eta f_0(\mathbf{y}) < f_1(\mathbf{y}) \\ \gamma(\mathbf{y}) & \text{if } f_1(\mathbf{y}) = \eta f_0(\mathbf{y}) \\ 0 & \text{if } f_1(\mathbf{y}) < \eta f_0(\mathbf{y}). \end{cases} \quad (5.6)$$

The inequality discussed next lays the groundwork for identifying the Neyman–Pearson test $\delta_{\text{NP}}(\alpha)$.

Lemma 5.3.1 *For any test $\delta : \mathbb{R}^k \rightarrow [0, 1]$ in \mathcal{D}^* , the inequality*

$$P_{\text{D}}(\delta^*) - P_{\text{D}}(\delta) \geq \eta (P_{\text{F}}(\delta^*) - P_{\text{F}}(\delta)) \quad (5.7)$$

holds where the randomized test $\delta^ : \mathbb{R}^k \rightarrow [0, 1]$ in \mathcal{D}^* is given by (5.7).*

If we select $\eta \geq 0$ and $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ so that δ^* satisfies the equality

$$P_{\text{F}}(\delta^*) = \alpha, \quad (5.8)$$

then the inequality (5.7) reads

$$P_{\text{D}}(\delta^*) - P_{\text{D}}(\delta) \geq \eta (\alpha - P_{\text{F}}(\delta)), \quad \delta \in \mathcal{D}^*. \quad (5.9)$$

For any test $\delta : \mathbb{R}^k \rightarrow [0, 1]$ in \mathcal{D}_{α}^* , we then conclude that

$$P_{\text{D}}(\delta^*) - P_{\text{D}}(\delta) \geq \eta (\alpha - P_{\text{F}}(\delta)) \geq 0 \quad (5.10)$$

since $P_{\text{F}}(\delta) \leq \alpha$. In other words,

$$P_{\text{D}}(\delta) \leq P_{\text{D}}(\delta^*), \quad \delta \in \mathcal{D}_{\alpha}^*$$

and the test δ^* solves the constrained problem NP_{α}^* .

We now show that the parameter $\eta \geq 0$ and the Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ can indeed be selected so that a test δ^* of the form (5.7) indeed satisfies (5.8).

Lemma 5.3.2 *For every α in $(0, 1]$ it is always possible to select $\eta \geq 0$ and a Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ in (5.6) so that (5.8) holds.*

Finally uniqueness is shown to hold in the following sense.

Lemma 5.3.3 *For every α in $(0, 1]$, if $\delta_{\text{NP}}(\alpha)$ is a Neyman–Pearson test (possibly in \mathcal{D}^*) of size α , then it necessarily holds that*

$$\mathbb{P}[\delta_{\text{NP}}(\alpha)(\mathbf{Y}) = \delta^*(\mathbf{Y}) | H = h] = 1, \quad h = 0, 1 \quad (5.11)$$

where the test δ^ is given by (5.6) with $\eta \geq 0$ and Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ selected so that (5.8) holds.*

5.4 Proofs

Throughout the discussion α is given in $(0, 1]$ and held fixed.

A proof of Lemma 5.3.1 Let $\delta : \mathbb{R}^k \rightarrow [0, 1]$ be an arbitrary test in \mathcal{D}^* . As discussed in Section 3.3 recall that

$$P_F(\delta) = \int_{\mathbb{R}^k} \delta(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) \quad \text{and} \quad P_F(\delta^*) = \int_{\mathbb{R}^k} \delta^*(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}),$$

while

$$P_D(\delta) = \int_{\mathbb{R}^k} \delta(\mathbf{y}) f_1(\mathbf{y}) dF(\mathbf{y}) \quad \text{and} \quad P_D(\delta^*) = \int_{\mathbb{R}^k} \delta^*(\mathbf{y}) f_1(\mathbf{y}) dF(\mathbf{y}).$$

It follows that

$$\begin{aligned} & P_D(\delta^*) - P_D(\delta) - \eta (P_F(\delta^*) - P_F(\delta)) \\ &= \int_{\mathbb{R}^k} (\delta^*(\mathbf{y}) - \delta(\mathbf{y})) f_1(\mathbf{y}) dF(\mathbf{y}) - \eta \int_{\mathbb{R}^k} (\delta^*(\mathbf{y}) - \delta(\mathbf{y})) f_0(\mathbf{y}) dF(\mathbf{y}) \\ &= \int_{\mathbb{R}^k} (\delta^*(\mathbf{y}) - \delta(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) dF(\mathbf{y}) \\ &= \int_{\mathbb{R}^k} P_\eta(\mathbf{y}) dF(\mathbf{y}) \end{aligned} \tag{5.12}$$

where we have set

$$P_\eta(\mathbf{y}) \equiv (\delta^*(\mathbf{y}) - \delta(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})), \quad \mathbf{y} \in \mathbb{R}^k.$$

Direct inspection shows that we always have

$$P_\eta(\mathbf{y}) \geq 0, \quad \mathbf{y} \in \mathbb{R}^k. \tag{5.13}$$

Obviously, we have $P_\eta(\mathbf{y}) = 0$ when $f_1(\mathbf{y}) = \eta f_0(\mathbf{y})$. When $\eta f_0(\mathbf{y}) < f_1(\mathbf{y})$, then

$$P_\eta(\mathbf{y}) = (1 - \delta(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) \geq 0,$$

while when $f_1(\mathbf{y}) < \eta f_0(\mathbf{y})$, then

$$P_\eta(\mathbf{y}) = -\delta(\mathbf{y}) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) \geq 0.$$

It is now plain from (5.12) and (5.13) that

$$P_D(\delta^*) - P_D(\delta) - \eta (P_F(\delta^*) - P_F(\delta)) \geq 0 \tag{5.14}$$

and the inequality (5.7) follows. ■

A proof of Lemma 5.3.2 Using the definition (5.6) of the randomized test δ^* , we note that

$$\begin{aligned}
& P_F(\delta^*) \\
&= \int_{\mathbb{R}^k} \delta^*(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) \\
&= \int_{\{\mathbf{y} \in \mathbb{R}^k: f_1(\mathbf{y}) = \eta f_0(\mathbf{y})\}} \gamma(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) + \int_{\{\mathbf{y} \in \mathbb{R}^k: f_1(\mathbf{y}) > \eta f_0(\mathbf{y})\}} f_0(\mathbf{y}) dF(\mathbf{y}) \\
&= \int_{\{\mathbf{y} \in \mathbb{R}^k: f_1(\mathbf{y}) = \eta f_0(\mathbf{y})\}} \gamma(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) + \mathbb{P}[f_1(\mathbf{Y}) > \eta f_0(\mathbf{Y}) | H = 0].
\end{aligned}$$

As we seek to satisfy (5.8), we need to select $\eta \geq 0$ and a Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ such that

$$\alpha - \mathbb{P}[f_1(\mathbf{Y}) > \eta f_0(\mathbf{Y}) | H = 0] = \int_{\{\mathbf{y} \in \mathbb{R}^k: f_1(\mathbf{y}) = \eta f_0(\mathbf{y})\}} \gamma(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}).$$

This last relation suggests introducing the quantity $\eta(\alpha)$ defined by

$$\eta(\alpha) = \inf \{ \eta \geq 0 : \mathbb{P}[f_1(\mathbf{Y}) > \eta f_0(\mathbf{Y}) | H = 0] < \alpha \}.$$

The definition of $\eta(\alpha)$ is well posed since $\eta \rightarrow \mathbb{P}[f_1(\mathbf{Y}) > \eta f_0(\mathbf{Y}) | H = 0]$ is non-increasing (and right-continuous) on $(0, \infty)$.

Two cases are possible: If

$$\mathbb{P}[f_1(\mathbf{Y}) > \eta(\alpha) f_0(\mathbf{Y}) | H = 0] < \alpha,$$

then take $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ to be constant, say

$$\gamma(\mathbf{y}) = \gamma(\alpha), \quad \mathbf{y} \in \mathbb{R}^k.$$

In that case, the constant $\gamma(\alpha)$ satisfies

$$\alpha - \mathbb{P}[f_1(\mathbf{Y}) > \eta(\alpha) f_0(\mathbf{Y}) | H = 0] = \gamma(\alpha) \int_{\{\mathbf{y} \in \mathbb{R}^k: f_1(\mathbf{y}) = \eta(\alpha) f_0(\mathbf{y})\}} f_0(\mathbf{y}) dF(\mathbf{y}),$$

whence

$$\gamma(\alpha) = \frac{\alpha - \mathbb{P}[f_1(\mathbf{Y}) > \eta(\alpha) f_0(\mathbf{Y}) | H = 0]}{\mathbb{P}[f_1(\mathbf{Y}) = \eta(\alpha) f_0(\mathbf{Y}) | H = 0]}. \quad (5.15)$$

If

$$\mathbb{P}[f_1(\mathbf{Y}) > \eta(\alpha)f_0(\mathbf{Y}) | H = 0] = \alpha,$$

then the mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ must be selected so that

$$\int_{\{\mathbf{y} \in \mathbb{R}^k : f_1(\mathbf{y}) = \eta f_0(\mathbf{y})\}} \gamma(\mathbf{y}) f_0(\mathbf{y}) dF(\mathbf{y}) = 0.$$

This can be achieved by taking the constant mapping given by

$$\gamma(\mathbf{y}) = 0, \quad \mathbf{y} \in \mathbb{R}^k.$$

■

A proof of Lemma 5.3.3 The test $\delta_{\text{NP}}(\alpha)$ being a Neyman–Pearson test of size α , the equality

$$P_{\text{D}}(\delta_{\text{NP}}(\alpha)) = P_{\text{D}}(\delta^*)$$

must hold where the test δ^* is given by (5.6) with $\eta > 0$ and Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$ selected so that (5.8) holds. This is a consequence of the fact that both $\delta_{\text{NP}}(\alpha)$ and δ^* solve the problem NP_α^* .

It then follows from (5.7) that

$$0 = P_{\text{D}}(\delta^*) - P_{\text{D}}(\delta_{\text{NP}}(\alpha)) \geq \eta(\alpha - P_{\text{F}}(\delta_{\text{NP}}(\alpha))) \geq 0 \quad (5.16)$$

since $P_{\text{F}}(\delta^*) = \alpha$ under the choice of $\eta > 0$ and the Borel mapping $\gamma : \mathbb{R}^k \rightarrow [0, 1]$, whence $P_{\text{F}}(\delta_{\text{NP}}(\alpha)) = \alpha$.

In other words, $P_{\text{D}}(\delta_{\text{NP}}(\alpha)) = P_{\text{D}}(\delta^*)$ and $P_{\text{F}}(\delta_{\text{NP}}(\alpha)) = P_{\text{F}}(\delta^*)$. Using these facts in the expression (5.12) (with the strategy $\delta_{\text{NP}}(\alpha)$) we find that

$$\begin{aligned} 0 &= P_{\text{D}}(\delta^*) - P_{\text{D}}(\delta_{\text{NP}}(\alpha)) - \eta(P_{\text{F}}(\delta^*) - P_{\text{F}}(\delta_{\text{NP}}(\alpha))) \\ &= \int_{\mathbb{R}^k} (\delta^*(\mathbf{y}) - \delta_{\text{NP}}(\alpha)(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) dF(\mathbf{y}) \end{aligned} \quad (5.17)$$

with

$$(\delta^*(\mathbf{y}) - \delta_{\text{NP}}(\alpha)(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) \geq 0, \quad \mathbf{y} \in \mathbb{R}^k$$

by virtue of (5.13). It immediately follows that

$$(\delta^*(\mathbf{y}) - \delta_{\text{NP}}(\alpha)(\mathbf{y})) (f_1(\mathbf{y}) - \eta f_0(\mathbf{y})) = 0 \quad F - a.e. \quad (5.18)$$

on \mathbb{R}^k . Therefore,

$$\delta_{\text{NP}}(\alpha)(\mathbf{y}) = \delta^*(\mathbf{y}) \quad F - a.e. \quad (5.19)$$

on $\{\mathbf{y} \in \mathbb{R}^k : f_1(\mathbf{y}) \neq \eta f_0(\mathbf{y})\}$. ■

5.5 Examples

The Gaussian case Consider again the situation discussed in Section 2.6 where the observation rv \mathbf{Y} is conditionally Gaussian given H , i.e.,

$$\begin{aligned} H_1 : \mathbf{Y} &\sim N(\mathbf{m}_1, \mathbf{R}) \\ H_0 : \mathbf{Y} &\sim N(\mathbf{m}_0, \mathbf{R}) \end{aligned}$$

where \mathbf{m}_1 and \mathbf{m}_0 are distinct elements in \mathbb{R}^k , and the $k \times k$ symmetric matrix \mathbf{R} is positive definite (thus invertible). From the discussion given in Section 2.6, it follows for each $\lambda > 0$ the test d_λ takes the form

$$d_\lambda(\mathbf{y}) = 0 \quad \text{iff} \quad \mathbf{y}'\mathbf{R}^{-1}\Delta\mathbf{m} > \phi(\lambda)$$

with $\Delta\mathbf{m}$ and $\phi(\lambda)$ given by (2.22) and (2.23), respectively. We also have

$$P_{\text{F}}(d_\lambda) = 1 - \Phi\left(\frac{\log \lambda + \frac{1}{2}d^2}{d}\right).$$

where d^2 is given by (2.28) – It is plain that the function $\lambda \rightarrow P_{\text{F}}(d_\lambda)$ is continuous on \mathbb{R}_+ with $\{P_{\text{F}}(d_\lambda), \lambda > 0\} = (0, 1)$. Given α in the unit interval $(0, 1)$, the value $\lambda(\alpha)$ is *uniquely* determined through the relation

$$1 - \alpha = \Phi\left(\frac{\log \lambda + \frac{1}{2}d^2}{d}\right).$$

This is equivalent to

$$\lambda(\alpha) = e^{d \cdot x_1 - \alpha - \frac{1}{2}d^2}.$$

where for t in $(0, 1)$, let x_t denote the only solution to the equation

$$\Phi(x) = t, \quad x \in \mathbb{R}.$$

Discontinuity with Bernoulli rvs The setting is that of Section 2.5 to which we refer the reader for the notation. We discuss only the case $a_1 < a_0$, and leave the case $a_0 < a_1$ as an exercise for the interested reader. We have shown that

$$P_F(d_\lambda) = \begin{cases} 1 & \text{if } 0 < \lambda \leq \frac{a_1}{a_0} \\ 1 - a_0 & \text{if } \frac{a_1}{a_0} < \lambda \leq \frac{1-a_1}{1-a_0} \\ 0 & \text{if } \frac{1-a_1}{1-a_0} < \lambda \end{cases} \quad (5.20)$$

as λ ranges over $(0, \infty)$.

Note that $\lambda \rightarrow P_F(d_\lambda)$ is left-continuous but not continuous with

$$\{P_F(d_\lambda), \lambda > 0\} = \{0, 1 - a_0, 1\}.$$

Discontinuity with Poisson rvs With $\mathcal{P}(m)$ denoting the Poisson pmf on \mathbb{N} with parameter $m > 0$, consider the following simple binary hypothesis testing problem

$$\begin{aligned} H_1 : & Y \sim \mathcal{P}(m_1) \\ H_0 : & Y \sim \mathcal{P}(m_0) \end{aligned}$$

where $m_1 \neq m_0$ in $(0, \infty)$. Thus,

$$\mathbb{P}[Y = k | H = h] = \frac{(m_h)^k}{k!} e^{-m_h}, \quad \begin{matrix} h = 0, 1 \\ k = 0, 1, \dots \end{matrix}$$

In this example, we take F to be the counting measure on \mathbb{N} , and for every $\lambda \geq 0$, the definition of d_λ reduces to

$$\begin{aligned} d_\lambda(k) = 0 & \text{ iff } \frac{(m_1)^k}{k!} e^{-m_1} < \lambda \frac{(m_0)^k}{k!} e^{-m_0} \\ & \text{ iff } \left(\frac{m_1}{m_0} \right)^k < \lambda e^{-(m_0 - m_1)} \end{aligned} \quad (5.21)$$

with $k = 0, 1, \dots$

If $m_0 < m_1$, then

$$\begin{aligned} d_\lambda(k) = 0 & \text{ iff } \frac{(m_1)^k}{k!} e^{-m_1} < \lambda \frac{(m_0)^k}{k!} e^{-m_0} \\ & \text{ iff } \left(\frac{m_1}{m_0} \right)^k < \lambda e^{-(m_0 - m_1)} \\ & \text{ iff } k < \eta(\lambda) \end{aligned} \quad (5.22)$$

with $k = 0, 1, \dots$, where

$$\eta(\lambda) = \frac{\log \lambda e^{-(m_0-m_1)}}{\log \left(\frac{m_1}{m_0} \right)}.$$

It follows that

$$\begin{aligned} P_{\text{F}}(d_\lambda) &= \mathbb{P}[d_\lambda(Y) = 1 | H = 0] \\ &= \mathbb{P}[Y \geq \eta(\lambda) | H = 0] \\ &= \sum_{k=0: \eta(\lambda) \leq k}^{\infty} \frac{(m_0)^k}{k!} e^{-m_0}. \end{aligned} \quad (5.23)$$

In this last expression only the integer ceiling $\lceil \eta(\lambda) \rceil$ of $\eta(\lambda)$ matters, where $\lceil \eta(\lambda) \rceil = \inf \{k \in \mathbb{N} : \eta(\lambda) \leq k\}$, whence

$$P_{\text{F}}(d_\lambda) = \sum_{k=\lceil \eta(\lambda) \rceil}^{\infty} \frac{(m_0)^k}{k!} e^{-m_0}.$$

As a result, the mapping $\lambda \rightarrow P_{\text{F}}(d_\lambda)$ is easily seen to be a *left-continuous piecewise constant* mapping with

$$P_{\text{F}}(d_\lambda) = P_{\text{F}}(d_{\lambda_n}), \quad \begin{array}{l} \lambda_n < \lambda \leq \lambda_{n+1} \\ n = 0, 1, \dots \end{array}$$

where $\{\lambda_n, n = 1, 2, \dots\}$ is a strictly monotone increasing sequence determined by the relation

$$n = \frac{\log \lambda_n e^{-(m_0-m_1)}}{\log \left(\frac{m_1}{m_0} \right)}, \quad n = 1, 2, \dots$$

or equivalently,

$$\lambda_n = \left(\frac{m_1}{m_0} \right)^n e^{-(m_1-m_0)}, \quad n = 1, 2, \dots$$

It is now plain that whenever α is chosen in $[0, 1]$ such that

for some integer $n = 0, 1, \dots$ then the requirement that $P_{\text{F}}(d_{\lambda(\alpha)}) = \alpha$ cannot be met. This difficulty is circumvented by enlarging \mathcal{D} with *randomized* policies; see Section ??.

5.6 Exercises**5.7 References**

Chapter 6

The receiver operating characteristics

In this chapter we investigate various properties of the mappings $\eta \rightarrow P_F(d_\eta)$ and $\eta \rightarrow P_F(d_\eta)$ as η ranges over \mathbb{R}_+ . This leads to defining the *receiver operating characteristic curve*, and to developing it into a handy operational tool to solve the various versions of the binary hypothesis problem discussed so far.

6.1 A basic limiting result

We start with a basic observation.

Lemma 6.1.1 *Assume Condition (A.1) to hold. For each $h = 0, 1$, the mapping*

$$\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow \mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h]$$

is monotone non-decreasing. Left (resp. right) limits exist at all points on $(0, \infty)$ (resp. $[0, \infty)$), with

$$\begin{aligned} & \lim_{\eta \downarrow \lambda} \mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] & (6.1) \\ = & \mathbb{P}[f_1(\mathbf{Y}) > \lambda f_0(\mathbf{Y}), f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P}[f_0(\mathbf{Y}) = 0 | H = h], \quad \lambda \geq 0 \end{aligned}$$

and

$$\begin{aligned} & \lim_{\eta \uparrow \lambda} \mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] \\ = & \mathbb{P}[f_1(\mathbf{Y}) \geq \lambda f_0(\mathbf{Y}) | H = h], \quad \lambda > 0. & (6.2) \end{aligned}$$

For future reference, for each $\eta \geq 0$ define the Borel subset $R(\eta)$ of \mathbb{R}^k by

$$R(\eta) \equiv \{\mathbf{y} \in \mathbb{R}^k : f_1(\mathbf{y}) \geq \eta f_0(\mathbf{y})\}. \quad (6.3)$$

Note that

$$\mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] = \mathbb{P}[\mathbf{Y} \in R(\eta) | H = h], \quad h = 0, 1, . \quad (6.4)$$

Proof. The asserted monotonicity property is a consequence of the inclusion $R(\eta_2) \subseteq R(\eta_1)$ holding whenever $\eta_1 < \eta_2$. The existence of left (resp. right) limits at all points on $(0, \infty)$ (resp. $[0, \infty)$) immediately follows.

Consider $h = 0, 1$. Fix $\lambda \geq 0$. By standard continuity facts from measure theory it follows that

$$\begin{aligned} \lim_{\eta \downarrow \lambda} \mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] &= \lim_{\eta \downarrow \lambda} \mathbb{P}[\mathbf{Y} \in R(\eta) | H = h] \\ &= \mathbb{P}[\cup_{\lambda < \eta} [\mathbf{Y} \in R(\eta)] | H = h] \\ &= \mathbb{P}[\mathbf{Y} \in \cup_{\eta > \lambda} R(\eta) | H = h], \end{aligned}$$

and we obtain (6.1) as we note that

$$\cup_{\eta > \lambda} R(\eta) = \left\{ \mathbf{y} \in \mathbb{R}^k : \begin{array}{l} f_1(\mathbf{y}) > \lambda f_0(\mathbf{y}), \\ f_0(\mathbf{y}) > 0 \end{array} \right\} \cup \{\mathbf{y} \in \mathbb{R}^k : f_0(\mathbf{y}) = 0\}.$$

In a very similar way, with $\lambda > 0$ we get

$$\begin{aligned} \lim_{\eta \uparrow \lambda} \mathbb{P}[f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] &= \lim_{\eta \uparrow \lambda} \mathbb{P}[\mathbf{Y} \in R(\eta) | H = h] \\ &= \mathbb{P}[\cap_{\eta < \lambda} [\mathbf{Y} \in R(\eta)] | H = h] \\ &= \mathbb{P}[\mathbf{Y} \in \cap_{\eta < \lambda} R(\eta) | H = h] \\ &= \mathbb{P}[f_1(\mathbf{Y}) \geq \lambda f_0(\mathbf{Y}) | H = h] \end{aligned}$$

because

$$\cap_{\eta < \lambda} R(\eta) = \{\mathbf{y} \in \mathbb{R}^k : f_1(\mathbf{y}) \geq \lambda f_0(\mathbf{y})\} = R(\lambda).$$

This establishes (6.2). ■

6.2 Continuity properties

We shall specialize Lemma 6.1.1 for $h = 0$ and $h = 1$.

Lemma 6.2.1 *The mapping $\eta \rightarrow P_{\mathbb{F}}(d_\eta)$ is monotone non-decreasing on \mathbb{R}_+ , with the left (resp. right) limit existing at all points in $(0, \infty)$ (resp. $[0, \infty)$). Under Conditions (A.1), it holds that In particular, it holds that*

$$\lim_{\eta \uparrow \lambda} P_{\mathbb{F}}(d_\eta) = \mathbb{P}[f_1(\mathbf{Y}) \geq \lambda f_0(\mathbf{Y}) | H = 0] = P_{\mathbb{F}}(d_\lambda), \quad \lambda > 0 \quad (6.5)$$

while

$$\lim_{\eta \downarrow \lambda} P_{\mathbb{F}}(d_\eta) = \mathbb{P}[f_1(\mathbf{Y}) > \lambda f_0(\mathbf{Y}) | H = 0], \quad \lambda \geq 0. \quad (6.6)$$

Proof. Applying (6.1) and (6.2) with $h = 0$ readily yields the desired conclusions as we recall that

$$\mathbb{P}[f_0(\mathbf{Y}) = 0 | H = 0] = 0 \quad (6.7)$$

under Condition (A.1). ■

Thus, the mapping $\eta \rightarrow P_{\mathbb{F}}(d_\eta)$ is left-continuous on $(0, \infty)$, but not necessarily right-continuous on $[0, \infty)$ as we note that

$$\lim_{\eta \downarrow \lambda} P_{\mathbb{F}}(d_\eta) = P_{\mathbb{F}}(d_\lambda) - \mathbb{P}[f_1(\mathbf{Y}) = \lambda f_0(\mathbf{Y}) | H = 0] \geq 0. \quad (6.8)$$

Lemma 6.2.2 *The mapping $\eta \rightarrow P_{\mathbb{D}}(d_\eta)$ is monotone non-decreasing on \mathbb{R}_+ , with the left (resp. right) limit existing at all points on $(0, \infty)$ (resp. $[0, \infty)$). Under Conditions (A.1) and (A.2), it holds that*

$$\lim_{\eta \uparrow \lambda} P_{\mathbb{D}}(d_\eta) = \mathbb{P}[f_1(\mathbf{Y}) \geq \lambda f_0(\mathbf{Y}) | H = 1] = P_{\mathbb{D}}(d_\lambda), \quad \lambda > 0 \quad (6.9)$$

while

$$\lim_{\eta \downarrow \lambda} P_{\mathbb{D}}(d_\eta) = \mathbb{P}[f_1(\mathbf{Y}) > \lambda f_0(\mathbf{Y}) | H = 1], \quad \lambda \geq 0. \quad (6.10)$$

Proof. As discussed in Section ??, under Conditions (A.1) and (A.2) it also holds that

$$\mathbb{P}[f_0(\mathbf{Y}) = 0|H = 1] = 0. \quad (6.11)$$

Applying Lemma 6.1.1 with $h = 1$ and using this last fact we readily get the result. ■

Again, the mapping $\eta \rightarrow P_D(d_\eta)$ is left-continuous on $(0, \infty)$ but not necessarily right-continuous on $[0, \infty)$ as we note that

$$\lim_{\eta \downarrow \lambda} P_D(d_\eta) = P_D(d_\lambda) - \mathbb{P}[f_1(\mathbf{Y}) = \lambda f_0(\mathbf{Y})|H = 1] \geq 0. \quad (6.12)$$

We close with the behavior at $\eta = 0$ and at $\eta = \infty$; proofs are available in Section 6.8.

Lemma 6.2.3 *Under Condition (A.1) we have*

$$\lim_{\eta \rightarrow 0} P_F(d_\eta) = \mathbb{P}[f_1(\mathbf{Y}) > 0|H = 0], \quad (6.13)$$

and if Condition (A.2) also holds, then

$$\lim_{\eta \rightarrow 0} P_D(d_\eta) = 1. \quad (6.14)$$

In principle, it is possible that $\mathbb{P}[f_1(\mathbf{Y}) > 0|H = 0] < 1$. However, with the notation (??) introduced in Section ??, we see that (6.13) becomes

$$\lim_{\eta \rightarrow 0} P_F(d_\eta) = 1 \quad (6.15)$$

if $B_0 \subseteq B_1$ since then $\mathbb{P}[f_1(\mathbf{Y}) > 0|H = 0] = 1$ (or equivalently, (6.7)) under Condition (A.1).

Note that $B_0 \subseteq B_1$ implies $B_0 = B_1$ under Conditions (A.1) and (A.2). In that case we have “continuity” at the origin because $P_F(d_0) = 1$ and $P_D(d_0) = 1$ (under the convention used for d_0 in Section ??).

Lemma 6.2.4 *Under Condition (A.1) we always have*

$$\lim_{\eta \rightarrow \infty} P_F(d_\eta) = 0, \quad (6.16)$$

and if Condition (A.2) also holds, then

$$\lim_{\eta \rightarrow \infty} P_D(d_\eta) = 0. \quad (6.17)$$

Note that $P_F(d_\infty) = 0$ and $P_D(d_\infty) = 0$ (under the convention used for d_∞ introduced Section ??), implying “continuity” at infinity.

6.3 The receiver operating characteristic (ROC)

A careful inspection of the solutions to the three formulations discussed so far shows that sometimes under mild assumptions, the test of interest takes the form

$$d_\eta(\mathbf{y}) = 0 \quad \text{iff} \quad f_1(\mathbf{y}) < \eta f_0(\mathbf{y}) \quad (6.18)$$

for some $\eta > 0$ – It is only the value of the threshold η that varies with the problem formulation. With the notation used earlier, we have

In the Bayesian formulation,

$$\eta_B = \frac{\Gamma_0(1-p)}{\Gamma_1 p}$$

In the minimax formulation,

$$\eta_m = \frac{\Gamma_0(1-p_m)}{\Gamma_1 p_m}$$

with p_m such that

$$V(p_m) = \max(V(p) : p \in [0, 1]).$$

When p_m is a point in $(0, 1)$ at which $V : [0, 1] \rightarrow \mathbb{R}$ is differentiable, then p_m can be characterized through the Minimax Equation.

In the Neyman–Pearson formulation,

$$\eta_{NP}(\alpha) = \lambda(\alpha).$$

with $\lambda(\alpha)$ satisfying the constraint (5.4).

In view of this, it seems natural to analyze in some details the performance of the tests (6.18). This is done by considering how their probabilities of detection and of false alarm, namely $P_F(d_\eta)$ and $P_D(d_\eta)$, vary in relation to each other as η ranges from $\eta = 0$ to $\eta = +\infty$. This is best understood by plotting the graph (Γ) of the detection probability against the corresponding probability of false alarm. Such a graph is analogous to a *phase portrait* for two-dimensional non-linear ODEs, and is called a *receiver operating characteristic* (ROC) curve. Its *parametric* representation is given by

$$\mathbb{R}_+ \rightarrow [0, 1] \times [0, 1] : \eta \rightarrow (P_F(d_\eta), (P_D(d_\eta))),$$

whence

$$(\Gamma) : \quad \{(P_F(d_\eta), (P_D(d_\eta))), \eta \geq 0\}.$$

This graph is *completely* determined by the probability distributions F_0 and F_1 of the observation rv \mathbf{Y} under the two hypotheses (through the densities f_0 and f_1 with respect to the underlying distribution F) and not by cost assignments or the prior probabilities. A typical ROC curve is drawn below.

6.4 Geometric properties of the ROC curve

The following geometric properties of the ROC curve are key to its operational usefulness.

Theorem 6.4.1 *Assume that Conditions (A.1) and (A.2) hold.*

(i): *Both mappings $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_F(d_\eta)$ and $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_D(d_\eta)$ are monotone non-increasing, with $\lim_{\eta \rightarrow 0} P_F(d_\eta) \equiv P_F^* \leq P_F(d_0) = 1$ and $\lim_{\eta \rightarrow 0} P_D(d_\eta) = P_D(d_\infty) = 1$, and $\lim_{\eta \uparrow \infty} P_F(d_\eta) = P_F(d_\infty) = 0$ and $\lim_{\eta \uparrow \infty} P_D(d_\eta) = P_D(d_\infty) = 0$.*

(ii): *If the right-derivative of $\eta \rightarrow P_F(d_\eta)$ exists at $\eta = \lambda$ for some $\lambda \geq 0$, then the right-derivative of $\eta \rightarrow P_D(d_\eta)$ also exists at $\eta = \lambda$, and the relation*

$$\left. \frac{d^+}{d\eta} P_D(d_\eta) \right|_{\eta=\lambda} = \lambda \cdot \left. \frac{d^+}{d\eta} P_F(d_\eta) \right|_{\eta=\lambda} \quad (6.19)$$

holds.

(iii): If the left-derivative of $\eta \rightarrow P_F(d_\eta)$ exists at $\eta = \lambda$ for some $\lambda > 0$, then the left-derivative of $\eta \rightarrow P_D(d_\eta)$ also exists at $\eta = \lambda$, and the relation

$$\left. \frac{d^-}{d\eta} P_D(d_\eta) \right|_{\eta=\lambda} = \lambda \cdot \left. \frac{d^-}{d\eta} P_F(d_\eta) \right|_{\eta=\lambda} \quad (6.20)$$

holds.

A proof of Theorem 6.4.1 is available in Section 6.7. It follows from this last result that whenever the mapping $\eta \rightarrow P_F(d_\eta)$ is *strictly decreasing* and *differentiable*, then the mapping $\eta \rightarrow P_D(d_\eta)$ is also *strictly decreasing* and *differentiable*, whence the curve (Γ) can be represented as the graph of a function $\Gamma : [0, 1] \rightarrow [0, 1] : P_F \rightarrow P_D = \Gamma(P_F)$, namely

$$P_D(d_\eta) = \Gamma(P_F(d_\eta)), \quad \eta \geq 0. \quad (6.21)$$

In such circumstance, Theorem 6.4.1 yields the the following information concerning this mapping. We consider only the case when $P_F^* = 1$ in Theorem 6.4.1; the situation when $P_F^* < 1$ can be handled in a similar way with details left to the interested reader.

Corollary 6.4.1 *Assume Conditions (A.1) and (A.2) to hold. Whenever the mapping $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_F(d_\eta)$ is differentiable and strictly decreasing, so is the mapping $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_D(d_\eta)$. In that case, there exists a mapping $\Gamma : [0, 1] \rightarrow [0, 1]$ satisfying (6.21) which is differentiable, strictly increasing and concave with*

$$\frac{d\Gamma}{dP_F}(P_F(d_\eta)) = \eta, \quad \eta \geq 0. \quad (6.22)$$

Proof. By Part (i) of Theorem 6.4.1 we see that the mapping $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_F(d_\eta)$ maps \mathbb{R}_+ onto $(0, 1]$. Being strictly decreasing and differentiable, this mapping admits an inverse, denoted $(0, 1] \rightarrow \mathbb{R}_+ : P_F \rightarrow \eta(P_F)$, with the property that

$$P_F(d_{\eta(P_F)}) = P_F, \quad P_F \in (0, 1].$$

By the Implicit Function Theorem this inverse mapping $P_F \rightarrow \eta(P_F)$ is differentiable; it is also strictly decreasing.

Define the mapping $\Gamma : [0, 1] \rightarrow [0, 1]$ by setting

$$\Gamma(P_F) \equiv P_D(d_{\eta(P_F)}), \quad P_F \in [0, 1] \quad (6.23)$$

with the understanding that $\Gamma(1) = 1$.

By Theorem 6.4.1 the differentiability of the mapping $\eta \rightarrow P_F(d_\eta)$ on \mathbb{R}_+ implies that of the mapping $\eta \rightarrow P_D(d_\eta)$ on \mathbb{R}_+ , and the mapping $P_F \rightarrow \Gamma(P_F)$ is therefore also differentiable. It is simple matter to check that this mapping is strictly decreasing on $[0, 1]$.

By the very definition of the function Γ , the identity (6.21) must hold. Differentiating both sides of (6.21) with respect of η we find

$$\begin{aligned} \frac{d}{d\eta} P_D(d_\eta) &= \frac{d}{d\eta} \Gamma(P_F(d_\eta)) \\ &= \frac{d\Gamma}{dP_F}(P_F(d_\eta)) \cdot \frac{d}{d\eta} P_F(d_\eta) \end{aligned} \quad (6.24)$$

as we use the Chain Rule. But Theorem 6.4.1 implies also that

$$\frac{d}{d\eta} P_D(d_\eta) = \eta \frac{d}{d\eta} P_F(d_\eta).$$

Combining these facts we conclude that

$$\eta \frac{d}{d\eta} P_F(d_\eta) = \frac{d\Gamma}{dP_F}(P_F(d_\eta)) \cdot \frac{d}{d\eta} P_F(d_\eta).$$

The mapping $\mathbb{R}_+ \rightarrow [0, 1] : \eta \rightarrow P_F(d_\eta)$ being assumed differentiable and strictly decreasing, we have $\frac{d}{d\eta} P_F(d_\eta) < 0$. Dividing by $\frac{d}{d\eta} P_F(d_\eta)$ we get (6.22). The other properties follow readily. ■

6.5 Operating the ROC

These results are most useful for operationally using the ROC curve:

For the Neyman–Pearson test of size α , consider the point on the ROC curve with abscissa α . It is determined by the threshold value $\eta(\alpha)$ with the property that $P_F(d_{\eta(\alpha)}) = \alpha$, and $d_{NP}(\alpha)$ is simply $d_{\eta(\alpha)}$. Note that $\eta(\alpha)$ is the *slope* of the tangent to the ROC curve at the point with abscissa α and the power $\beta(\alpha)$ of the test is simply the ordinate of that point.

For the Bayesian problem, η is determined by the cost assignment and the prior distribution of the rv H . The values of $P_D(d_\eta)$ and $P_F(d_\eta)$ can be easily determined by finding the point on the ROC where the tangent has slope η .

The Minimax Equation takes the form

$$C(1, 1) - C(0, 0) + \Gamma_1 P_M(d_\eta) - \Gamma_0 P_F(d_\eta) = 0,$$

or equivalently

$$C(1, 1) - C(0, 0) + \Gamma_1 = \Gamma_1 P_D(d_\eta) + \Gamma_0 P_F(d_\eta).$$

This shows that the minimax rule d_m^* is obtained as follows. Consider the straight line (L) in the (P_F, P_D) -plane with equation

$$(L) \quad C(1, 1) - C(0, 0) + \Gamma_1 = \Gamma_1 P_D + \Gamma_0 P_F.$$

Let (P_F^*, P_D^*) be the point of intersection of the straight line (L) with the ROC curve (Γ), and let η^* be the corresponding threshold value, i.e., $P_F^* = P_F(d_{\eta^*})$ and $P_D^* = P_D(d_{\eta^*})$. It is now clear that $d_m^* = d_{\eta^*}$.

6.6 Examples

Building on material developed earlier we now discuss the ROC in the Gaussian and Bernoulli cases, respectively.

The Gaussian case The setting is that of Section 2.4 to which we refer the reader for the notation. As shown there, for any $\eta > 0$ we have

$$P_F(d_\eta) = 1 - \Phi\left(\frac{\log \eta + \frac{1}{2}d^2}{d}\right)$$

and

$$P_D(d_\eta) = 1 - \Phi\left(\frac{\log \eta - \frac{1}{2}d^2}{d}\right).$$

To find the ROC curve, note that

$$d\Phi^{-1}(1 - P_F(d_\eta)) - \frac{d^2}{2} = \log \eta,$$

while we must have

$$d\Phi^{-1}(1 - P_D(d_\eta)) + \frac{d^2}{2} = \log \eta,$$

whence

$$d\Phi^{-1}(1 - P_F(d_\eta)) - \frac{d^2}{2} = d\Phi^{-1}(1 - P_D(d_\eta)) + \frac{d^2}{2}.$$

It follows that

$$\Phi^{-1}(1 - P_D(d_\eta)) = \Phi^{-1}(1 - P_F(d_\eta)) - d$$

so that

$$1 - P_D(d_\eta) = \Phi(\Phi^{-1}(1 - P_F(d_\eta)) - d)$$

This shows that here the mapping $\Gamma : [0, 1] \rightarrow [0, 1]$ is well defined and given by

$$P_D = 1 - \Phi(\Phi^{-1}(1 - P_F) - d), \quad P_F \in [0, 1].$$

The Bernoulli case The setting is that of Section 2.5 to which we refer the reader for the notation. We discuss only the case $a_1 < a_0$, and leave the case $a_0 < a_1$ as an exercise for the interested reader. It was shown that

$$P_F(d_\eta) = \begin{cases} 1 & \text{if } 0 < \eta \leq \frac{a_1}{a_0} \\ 1 - a_0 & \text{if } \frac{a_1}{a_0} < \eta \leq \frac{1-a_1}{1-a_0} \\ 0 & \text{if } \frac{1-a_1}{1-a_0} < \eta \end{cases} \quad (6.25)$$

and

$$P_D(d_\eta) = \begin{cases} 1 & \text{if } 0 < \eta \leq \frac{a_1}{a_0} \\ 1 - a_1 & \text{if } \frac{a_1}{a_0} < \eta \leq \frac{1-a_1}{1-a_0} \\ 0 & \text{if } \frac{1-a_1}{1-a_0} < \eta. \end{cases} \quad (6.26)$$

Therefore,

$$\{(P_F(d_\eta), P_D(d_\eta)), \eta \geq 0\} = \{(0, 0), (1 - a_0, 1 - a_1), (1, 1)\}.$$

Strictly speaking, in this case the ROC is not a “curve” as it comprises only three points. However, the points on the two segments $[(0, 0), (1 - a_0, 1 - a_1)]$ and $[(1 - a_0, 1 - a_1), (1, 1)]$ are achievable through randomization.

6.7 A proof of Theorem 6.4.1

We start with some facts that prove useful in discussing Theorem 6.4.1: With $\lambda > 0$, recall the set $R(\lambda)$ defined in Section 6.2 as

$$R(\lambda) \equiv \{\mathbf{y} \in \mathbb{R}^k : f_1(\mathbf{y}) \geq \lambda f_0(\mathbf{y})\}.$$

Noting that

$$d_\lambda(\mathbf{y}) = 1 \quad \text{iff} \quad \mathbf{y} \in R(\lambda),$$

it is plain that

$$\begin{aligned} P_F(d_\lambda) &= \mathbb{P}[d_\lambda(\mathbf{Y}) = 1 | H = 0] \\ &= \mathbb{P}[\mathbf{Y} \in R(\lambda) | H = 0] = \int_{R(\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) \end{aligned}$$

and

$$P_D(d_\lambda) = \mathbb{P}[d_\lambda(\mathbf{Y}) = 1 | H = 1] = \int_{R(\lambda)} f_1(\mathbf{y}) dF(\mathbf{y}).$$

For each $\Delta\lambda > 0$, easy algebra now leads to

$$\begin{aligned} P_F(d_{\lambda+\Delta\lambda}) - P_F(d_\lambda) &= \int_{R(\lambda+\Delta\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) - \int_{R(\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) \\ &= - \int_{R_+(\lambda; \Delta\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) \end{aligned}$$

where

$$R_+(\lambda; \Delta\lambda) \equiv \{\mathbf{y} \in \mathbb{R}^k : \lambda f_0(\mathbf{y}) \leq f_1(\mathbf{y}) < (\lambda + \Delta\lambda) f_0(\mathbf{y})\}.$$

Similarly, we have

$$\begin{aligned} P_F(d_{\lambda-\Delta\lambda}) - P_F(d_\lambda) &= \int_{R(\lambda-\Delta\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) - \int_{R(\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) \\ &= \int_{R_-(\lambda; \Delta\lambda)} f_0(\mathbf{y}) dF(\mathbf{y}) \end{aligned}$$

where

$$R_-(\lambda; \Delta\lambda) \equiv \{\mathbf{y} \in \mathbb{R}^k : (\lambda - \Delta\lambda) f_0(\mathbf{y}) \leq f_1(\mathbf{y}) < \lambda f_0(\mathbf{y})\}.$$

We can now turn to the proof of Theorem 6.4.1: Part (i) is already covered in Section 6.2. We shall discuss only Part (ii) as Part (iii) can be established by similar arguments *mutatis mutandi*. This is left to the interested reader.

Fix $\eta \geq 0$. With $\Delta\eta > 0$, the very definition of $R_+(\eta; \Delta\eta)$ implies the inequalities

$$\eta \int_{R_+(\eta; \Delta\eta)} f_0(\mathbf{y}) dF(\mathbf{y}) \leq \int_{R_+(\eta; \Delta\eta)} f_1(\mathbf{y}) dF(\mathbf{y})$$

and

$$\int_{R_+(\eta; \Delta\eta)} f_1(\mathbf{y}) dF(\mathbf{y}) \leq (\eta + \Delta\eta) \int_{R_+(\eta; \Delta\eta)} f_0(\mathbf{y}) dF(\mathbf{y}).$$

It then follows that

$$(\eta + \Delta\eta) \cdot \frac{P_F(d_{\eta+\Delta\eta}) - P_F(d_\eta)}{\Delta\eta} \leq \frac{P_D(d_{\eta+\Delta\eta}) - P_D(d_\eta)}{\Delta\eta}$$

and

$$\frac{P_D(d_{\eta+\Delta\eta}) - P_D(d_\eta)}{\Delta\eta} \leq \eta \cdot \frac{P_F(d_{\eta+\Delta\eta}) - P_F(d_\eta)}{\Delta\eta}.$$

If the right-derivative of $\eta \rightarrow P_F(d_\eta)$ exists, then

$$\frac{d^+}{d\eta} P_F(d_\eta) = \lim_{\Delta\eta \downarrow 0} \frac{P_F(d_{\eta+\Delta\eta}) - P_F(d_\eta)}{\Delta\eta}$$

and an easy sandwich argument shows that the limit

$$\lim_{\Delta\eta \downarrow 0} \frac{P_D(d_{\eta+\Delta\eta}) - P_D(d_\eta)}{\Delta\eta}$$

also exists. Therefore, the right-derivative of $\eta \rightarrow P_D(d_\eta)$ also exists and is given by

$$\frac{d^+}{d\eta} P_D(d_\eta) = \eta \cdot \frac{d^+}{d\eta} P_F(d_\eta).$$

■

6.8 Proofs of Lemmas 6.2.3 and 6.2.4

A proof of Lemma 6.2.3 Applying (6.1) (with $\lambda = 0$) yields

$$\begin{aligned} & \lim_{\eta \downarrow 0} \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] \\ &= \mathbb{P} [f_1(\mathbf{Y}) > 0, f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P} [f_0(\mathbf{Y}) = 0 | H = h], \quad h = 0, 1. \end{aligned} \quad (6.27)$$

Under Condition **(A.1)**, the validity of (6.7) implies (6.13) as we use (6.27) with $h = 0$. Under the additional Condition **(A.2)**, (6.11) and (6.27), this time with $h = 1$, lead to

$$\lim_{\eta \downarrow 0} \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = 1] = \mathbb{P} [f_1(\mathbf{Y}) > 0 | H = 1]$$

and (6.14) readily follows. ■

A proof of Lemma 6.2.4 Fix $h = 0, 1$. For each $\eta > 0$, we note that

$$\begin{aligned} & \mathbb{P} [d_\eta(\mathbf{Y}) = 1 | H = h] \\ &= \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}) | H = h] \\ &= \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}), f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P} [f_1(\mathbf{Y}) \geq 0, f_0(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}), f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P} [f_0(\mathbf{Y}) = 0 | H = h]. \end{aligned}$$

The usual monotonicity argument yields

$$\begin{aligned} & \lim_{\eta \rightarrow \infty} \mathbb{P} [d_\eta(\mathbf{Y}) = 1 | H = h] \\ &= \lim_{\eta \rightarrow \infty} \mathbb{P} [f_1(\mathbf{Y}) \geq \eta f_0(\mathbf{Y}), f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P} [f_0(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P} [f_1(\mathbf{Y}) = \infty, f_0(\mathbf{Y}) > 0 | H = h] + \mathbb{P} [f_0(\mathbf{Y}) = 0 | H = h] \\ &= \mathbb{P} [f_0(\mathbf{Y}) = 0 | H = h]. \end{aligned} \quad (6.28)$$

First use (6.28) with $h = 0$: Under Condition **(A.1)** the consequence (6.7) implies $\lim_{\eta \rightarrow \infty} \mathbb{P} [d_\eta(\mathbf{Y}) = 1 | H = 0] = 0$. With $h = 1$, under the additional Condition **(A.2)** we also get $\lim_{\eta \rightarrow \infty} \mathbb{P} [d_\eta(\mathbf{Y}) = 1 | H = 1] = 0$. This completes the proof of Lemma 6.2.4. ■

6.9 Exercises

6.10 References

Chapter 7

The M -ary hypothesis testing problem

As we shall see in this chapter and in the next one, the simple binary hypothesis testing problem of Chapter 7 admits several important generalizations. The version discussed here assumes that there are more than two hypotheses; it is of particular relevance to the design of optimal receivers in digital modulation. In this formulation, nature assumes M distinct states, labeled $0, 1, \dots, M - 1$, and is now encoded in a rv H which take values in the discrete set $\{0, 1, \dots, M - 1\}$. A decision has to be made as to which of these M hypotheses is the correct one on the basis of an observation \mathbf{Y} which is statistically related to H .

7.1 Motivating examples

Digital communications

Manufacturing

7.2 The probabilistic model

To formulate the M -ary hypothesis testing problem we proceed very much as in Chapter 2: With positive integer $M \geq 2$, we are given M distinct probability distribution functions F_0, \dots, F_{M-1} on \mathbb{R}^k , and a pmf $\mathbf{p} = (p_0, \dots, p_{M-1})$ on

$\{0, 1, \dots, M-1\}$ with

$$0 \leq p_m \leq 1, \quad m = 0, 1, \dots, M-1 \quad \text{and} \quad \sum_{m=0}^{M-1} p_m = 1.$$

The situation is summarized by

$$H_m : \mathbf{Y} \sim F_m, \quad m = 0, 1, \dots, M-1. \quad (7.1)$$

We construct a sample space Ω equipped with a σ -field of events \mathcal{F} , and rvs H and \mathbf{Y} defined on it and taking values in $\{0, 1, \dots, M-1\}$ and \mathbb{R}^k , respectively. Now the probability distribution functions F_0, \dots, F_{M-1} have the interpretation that

$$F_m(\mathbf{y}) = \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = m], \quad \mathbf{y} \in \mathbb{R}^k, \quad m = 0, 1, \dots, M-1.$$

The probability distribution of the rv H is specified by the pmf $\mathbf{p} = (p_0, \dots, p_{M-1})$ with

$$p_m = \mathbb{P}[H = m], \quad m = 0, 1, \dots, M-1.$$

Again, the conditional probability distributions of the observations given the hypothesis *and* the probability distribution of H completely specify the *joint* distribution of the rvs H and \mathbf{Y} : Indeed, for each $m = 0, 1, \dots, M-1$,

$$\begin{aligned} \mathbb{P}[\mathbf{Y} \leq \mathbf{y}, H = m] &= \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = m] \mathbb{P}[H = m] \\ &= p_m F_m(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned} \quad (7.2)$$

The *unconditional* probability distribution function of the rv \mathbf{Y} is easily determined to be

$$\mathbb{P}[\mathbf{Y} \leq \mathbf{y}] = \sum_{m=0}^{M-1} \mathbb{P}[\mathbf{Y} \leq \mathbf{y}, H = m], \quad \mathbf{y} \in \mathbb{R}^k$$

by the law of total probabilities, whence

$$\begin{aligned} G(\mathbf{y}) &\equiv \mathbb{P}[\mathbf{Y} \leq \mathbf{y}] \\ &= \sum_{m=0}^{M-1} p_m F_m(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned} \quad (7.3)$$

During the discussion, several assumptions will be enforced on the probability distributions F_0, \dots, F_{M-1} . The most common assumption is denoted by **(A.3)** for sake of convenience; it parallels Condition **(A.1)** made in the binary case:

Condition (A.3): The probability distributions F_0, \dots, F_{M-1} on \mathbb{R}^k are *absolutely continuous* with respect to some distribution F on \mathbb{R}^k .

This is equivalent to saying that for each $m = 0, 1, \dots, M-1$, there exists a Borel mapping $f_m : \mathbb{R}^k \rightarrow \mathbb{R}_+$ such that

$$F_m(\mathbf{y}) = \int_{-\infty}^{\mathbf{y}} f_m(\boldsymbol{\eta}) dF(\boldsymbol{\eta}), \quad \mathbf{y} \in \mathbb{R}^k, \quad m = 0, 1, \dots, M-1. \quad (7.4)$$

We refer to these Borel mappings as probability density functions with respect to F .

This condition is hardly constraining since we can always take F to be the average of the M distributions F_0, \dots, F_{M-1} . i.e.,

$$F(\mathbf{y}) \equiv \frac{1}{M} \left(\sum_{m=0}^{M-1} F_m(\mathbf{y}) \right), \quad \mathbf{y} \in \mathbb{R}^k \quad (7.5)$$

However, in most applications F is either Lebesgue measure on \mathbb{R}^k or a counting measure on some countable subset of \mathbb{R}^k .

Under Condition (A.3), the unconditional probability distribution function $G : \mathbb{R}^k \rightarrow [0, 1]$ of the rv \mathbf{Y} is automatically absolutely continuous with respect to the distribution F on \mathbb{R}^k : Indeed, we see from (7.4) and (7.3) that

$$\begin{aligned} G(\mathbf{y}) &= \sum_{m=0}^{M-1} p_m \int_{-\infty}^{\mathbf{y}} f_m(\boldsymbol{\eta}) dF(\boldsymbol{\eta}) \\ &= \int_{-\infty}^{\mathbf{y}} \left(\sum_{m=0}^{M-1} p_m f_m(\boldsymbol{\eta}) \right) dF(\boldsymbol{\eta}) \\ &= \int_{-\infty}^{\mathbf{y}} g(\boldsymbol{\eta}) dF(\boldsymbol{\eta}). \quad \mathbf{y} \in \mathbb{R}^k \end{aligned} \quad (7.6)$$

with Borel mapping $g : \mathbb{R}^k \rightarrow \mathbb{R}_+$ given by

$$g(\mathbf{y}) = \sum_{m=0}^{M-1} p_m f_m(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \quad (7.7)$$

In other words, the unconditional probability distribution function $G : \mathbb{R}^k \rightarrow [0, 1]$ admits $g : \mathbb{R}^k \rightarrow \mathbb{R}_+$ as probability density function with respect to F .

7.3 Admissible tests

We begin with a formal definition of an admissible test in the context M -ary hypothesis testing.

An *admissible* decision rule or test is any *Borel* mapping $d : \mathbb{R}^k \rightarrow \{0, 1, \dots, M-1\}$. The collection of all such admissible rules is still denoted by \mathcal{D} .

Again the measurability requirement is imposed to guarantee that the mapping $d(Y) : \Omega \rightarrow \{0, 1, \dots, M-1\} : \omega \rightarrow d(Y(\omega))$ is indeed a rv, i.e., $\{\omega \in \Omega : d(Y(\omega)) = m\}$ is an event in \mathcal{F} for all $m = 0, 1, \dots, M-1$.

A collection $\{C_0, \dots, C_{M-1}\}$ of M subsets of \mathbb{R}^k forms an *M -ary Borel partition* of \mathbb{R}^k if

- (i) For each $m = 0, 1, \dots, M-1$, the set C_m is a Borel subset of \mathbb{R}^k ;
- (ii) We have

$$C_m \cap C_k = \emptyset, \quad m \neq k \\ m, k = 0, 1, \dots, M-1$$

and

- (iii) The condition

$$\cup_{m=0}^{M-1} C_m = \mathbb{R}^k$$

holds.

The collection of all M -ary Borel partitions of \mathbb{R}^k is denoted $\mathcal{P}_M(\mathbb{R}^k)$. Lemma 1.5.1 has the following analog in the context M -ary hypothesis testing.

Lemma 7.3.1 *The set \mathcal{D} of admissible decision rules is in one-to-one correspondence with $\mathcal{P}_M(\mathbb{R}^k)$.*

Proof. For every test d in \mathcal{D} , the Borel sets defined by

$$C_m(d) \equiv \{\mathbf{y} \in \mathbb{R}^k : d(\mathbf{y}) = m\}, \quad m = 0, 1, \dots, M-1$$

are determined by d , and obviously form an M -ary Borel partition of \mathbb{R}^k .

Conversely, consider an M -ary Borel measurable partition $\{C_0, \dots, C_{M-1}\}$ in $\mathcal{P}_M(\mathbb{R}^k)$. With this partition we can associate the mapping $d_{C_0, \dots, C_{M-1}} : \mathbb{R}^k \rightarrow \{0, \dots, M-1\}$ given by

$$d_{C_0, \dots, C_{M-1}}(\mathbf{y}) = m, \quad \mathbf{y} \in C_m, \\ m = 0, \dots, M-1.$$

By construction this mapping is an admissible test in \mathcal{D} as we note that

$$C_m(d_{C_0, \dots, C_{M-1}}) = C_m, \quad m = 0, \dots, M-1$$

by the fact that the collection $\{C_0, \dots, C_{M-1}\}$ is a partition of \mathbb{R}^k . ■

7.4 The Bayesian formulation

The probabilistic model The Bayesian formulation assumes *knowledge* of the prior distribution $\mathbf{p} = (p_0, \dots, p_{M-1})$ of the rv H , and of the conditional distributions F_0, \dots, F_{M-1} of the rv \mathbf{Y} given H .

The optimization problem A cost is incurred for making decisions. This is captured through the mapping $C : \{0, 1, \dots, M-1\} \times \{0, 1, \dots, M-1\} \rightarrow \mathbb{R}$ with the interpretation that

$$C(m, \ell) = \begin{array}{l} \text{Cost incurred for deciding } \ell \\ \text{when } H = m \end{array}, \quad \ell, m = 0, 1, \dots, M-1.$$

Using the admissible rule d in \mathcal{D} incurs a cost $C(H, d(Y))$, but as for the binary hypothesis testing problem, the value of the cost $C(H, d(\mathbf{Y}))$ is not available, and attention focuses on the *expected cost* $J : \mathcal{D} \rightarrow \mathbb{R}$ given by

$$J(d) \equiv \mathbb{E}[C(H, d(\mathbf{Y}))], \quad d \in \mathcal{D}.$$

Here as well, the *Bayesian Problem* \mathcal{P}_B is the minimization problem

$$\mathcal{P}_B : \quad \text{Minimize } J(d) \text{ over } d \text{ in } \mathcal{D}.$$

Its solution is any test d^* in \mathcal{D} such that

$$J(d^*) \leq J(d), \quad d \in \mathcal{D}. \tag{7.8}$$

Any test d^* in \mathcal{D} which satisfies (7.8) is called a Bayesian test, and the value

$$J(d^*) = \inf_{d \in \mathcal{D}} J(d) = \min_{d \in \mathcal{D}} J(d) \tag{7.9}$$

is known as the *Bayesian cost*.

The Bayesian test For each $\ell = 0, 1, \dots, M - 1$, we define the Borel mapping $h_\ell : \mathbb{R}^k \rightarrow \mathbb{R}$ given by

$$h_\ell(\mathbf{y}) \equiv \sum_{m=0, m \neq \ell}^{M-1} p_m (C(m, \ell) - C(m, m)) f_m(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \quad (7.10)$$

The next result identifies the Bayesian test; its proof parallels that given in Section 2.2 for the binary case but with some important differences. Details are given in Section 7.5.

Theorem 7.4.1 *Assume the absolute continuity condition (A.1) to hold. The test $d^* : \mathbb{R}^k \rightarrow \{0, 1, \dots, M - 1\}$ given by*

$$d^*(\mathbf{y}) = \arg \min (\ell = 0, \dots, M - 1 : h_\ell(\mathbf{y})), \quad \mathbf{y} \in \mathbb{R}^k \quad (7.11)$$

(with a lexicographic tiebreaker in the event of ties) is admissible and solves the Bayesian Problem \mathcal{P}_B .

7.5 A proof of Theorem 7.4.1

A reduction step Fix a test d in \mathcal{D} . The decomposition

$$\mathbf{1} [d(\mathbf{Y}) = H] + \mathbf{1} [d(\mathbf{Y}) \neq H] = 1$$

holds so that

$$\begin{aligned} C(H, d(\mathbf{Y})) &= \mathbf{1} [d(\mathbf{Y}) = H] C(H, H) + \mathbf{1} [d(\mathbf{Y}) \neq H] C(H, d(\mathbf{Y})) \\ &= C(H, H) + (C(H, d(\mathbf{Y})) - C(H, H)) \mathbf{1} [d(\mathbf{Y}) \neq H]. \end{aligned} \quad (7.12)$$

Defining the auxiliary expected cost function $\hat{J} : \mathcal{D} \rightarrow \mathbb{R}$ to be

$$\hat{J}(d) = \mathbb{E} [(C(H, d(\mathbf{Y})) - C(H, H)) \mathbf{1} [d(\mathbf{Y}) \neq H]], \quad d \in \mathcal{D} \quad (7.13)$$

we again readily conclude that

$$J(d) = E[C(H, H)] + \hat{J}(d), \quad d \in \mathcal{D}. \quad (7.14)$$

Therefore, solving \mathcal{P}_B is equivalent to solving the auxiliary problem $\hat{\mathcal{P}}_B$ where

$$\hat{\mathcal{P}}_B : \quad \text{Minimize } \hat{J}(d) \text{ over } d \text{ in } \mathcal{D}.$$

Preparatory computations Fix d in \mathcal{D} . From (7.14) we note that

$$\begin{aligned}
\widehat{J}(d) &= \mathbb{E} [(C(H, d(\mathbf{Y})) - C(H, H)) \mathbf{1} [d(\mathbf{Y}) \neq H]] \\
&= \sum_{m=0}^{M-1} \left(\sum_{\ell=0, \ell \neq m}^{M-1} \mathbb{E} [(C(m, \ell) - C(m, m)) \mathbf{1} [H = m, d(\mathbf{Y}) = \ell]] \right) \\
&= \sum_{m=0}^{M-1} \left(\sum_{\ell=0, \ell \neq m}^{M-1} (C(m, \ell) - C(m, m)) \mathbb{P} [H = m, d(\mathbf{Y}) = \ell] \right) \\
&= \sum_{m=0}^{M-1} \left(\sum_{\ell=0, \ell \neq m}^{M-1} (C(m, \ell) - C(m, m)) \mathbb{P} [d(\mathbf{Y}) = \ell | H = m] p_m \right) \\
&= \sum_{m=0}^{M-1} \left(\sum_{\ell=0, \ell \neq m}^{M-1} p_m (C(m, \ell) - C(m, m)) \int_{C_\ell(d)} f_m(\mathbf{y}) dF(\mathbf{y}) \right).
\end{aligned}$$

Interchanging the order of summation we get

$$\begin{aligned}
\widehat{J}(d) &= \sum_{\ell=0}^{M-1} \left(\sum_{m=0, m \neq \ell}^{M-1} p_m (C(m, \ell) - C(m, m)) \int_{C_\ell(d)} f_m(\mathbf{y}) dF(\mathbf{y}) \right) \\
&= \sum_{\ell=0}^{M-1} \int_{C_\ell(d)} \left(\sum_{m=0, m \neq \ell}^{M-1} p_m (C(m, \ell) - C(m, m)) f_m(\mathbf{y}) \right) dF(\mathbf{y}) \\
&= \sum_{\ell=0}^{M-1} \int_{C_\ell(d)} h_\ell(\mathbf{y}) dF(\mathbf{y}) \tag{7.15}
\end{aligned}$$

where for each $\ell = 0, 1, \dots, M-1$, the mapping $h_\ell : \mathbb{R}^k \rightarrow \mathbb{R}$ is given by (7.10).

Solving $\widehat{\mathcal{P}}_B$ Define the mapping $h : \mathbb{R}^k \rightarrow \mathbb{R}$ given by

$$h(\mathbf{y}) \equiv \min_{\ell=0,1,\dots,M-1} h_\ell(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k. \tag{7.16}$$

where for each $\ell = 0, 1, \dots, M-1$, the mapping $h_\ell : \mathbb{R}^k \rightarrow \mathbb{R}$ is given by (7.10). The following facts

$$h(\mathbf{y}) = h_{d^*}(\mathbf{y})(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k \tag{7.17}$$

and

$$h(\mathbf{y}) \leq h_\ell(\mathbf{y}), \quad \begin{array}{l} \ell = 0, 1, \dots, M-1, \\ \mathbf{y} \in \mathbb{R}^k \end{array} \tag{7.18}$$

are simple consequences of the definition (7.11) of d^* .

Pick an arbitrary test in \mathcal{D} . Using the expression (7.15) we get

$$\begin{aligned}
& J(d) - J(d^*) \\
&= \widehat{J}(d) - \widehat{J}(d^*) \\
&= \sum_{\ell=0}^{M-1} \int_{C_\ell(d)} h_\ell(\mathbf{y}) dF(\mathbf{y}) - \sum_{\ell=0}^{M-1} \int_{C_\ell(d^*)} h_\ell(\mathbf{y}) dF(\mathbf{y}) \\
&\geq \sum_{\ell=0}^{M-1} \int_{C_\ell(d)} h(\mathbf{y}) dF(\mathbf{y}) - \sum_{\ell=0}^{M-1} \int_{C_\ell(d^*)} h_\ell(\mathbf{y}) dF(\mathbf{y}) \\
&= \sum_{\ell=0}^{M-1} \int_{C_\ell(d)} h(\mathbf{y}) dF(\mathbf{y}) - \sum_{\ell=0}^{M-1} \int_{C_\ell(d^*)} h(\mathbf{y}) dF(\mathbf{y}) \\
&= \int_{\mathbb{R}^k} h(\mathbf{y}) dF(\mathbf{y}) - \int_{\mathbb{R}^k} h(\mathbf{y}) dF(\mathbf{y}) \\
&= 0.
\end{aligned} \tag{7.19}$$

The inequality above and the equality that follows are consequences of (7.17) and (7.18), respectively. The last two steps used the fact that the collection $\{C_0(d), \dots, C_{M-1}(d)\}$ (resp. $\{C_0(d^*), \dots, C_{M-1}(d^*)\}$) forms an M -ary Borel partition of \mathbb{R}^k . In particular, this observation implies

$$\sum_{\ell=0}^{M-1} \int_{C_\ell(d)} h(\mathbf{y}) dF(\mathbf{y}) = \int_{\mathbb{R}^k} h(\mathbf{y}) dF(\mathbf{y})$$

and

$$\sum_{\ell=0}^{M-1} \int_{C_\ell(d^*)} h(\mathbf{y}) dF(\mathbf{y}) = \int_{\mathbb{R}^k} h(\mathbf{y}) dF(\mathbf{y}).$$

From (7.19) it follows that

$$J(d) - J(d^*) \geq 0, \quad d \in \mathcal{D},$$

and the optimality of d^* is now established. ■

7.6 The probability of error criterion

When C takes the form

$$C(m, k) = \mathbf{1}[\ell \neq m], \quad m, \ell = 0, 1, \dots, M-1,$$

the expected cost reduces to the *probability of error* criterion given by

$$P_E(d) = \mathbb{P}[d(\mathbf{Y}) \neq H], \quad d \in \mathcal{D}.$$

The Bayesian test $d^* : \mathbb{R}^k \rightarrow \{0, 1, \dots, M-1\}$ given by (7.11) now takes the following form: For each $\ell = 0, \dots, M-1$, the mapping $h_\ell : \mathbb{R}^k \rightarrow \mathbb{R}$ is now given by

$$h_\ell(\mathbf{y}) = \sum_{m=0, m \neq \ell}^{M-1} p_m f_m(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k.$$

But the probability distribution function $G : \mathbb{R}^k \rightarrow [0, 1]$ of the observation rv \mathbf{Y} is given by (7.3), and under condition **(A.1)** it has probability density function $g : \mathbb{R}^k \rightarrow \mathbb{R}_+$ given by (7.7). Therefore, for each $\ell = 0, \dots, M-1$, we have

$$h_\ell(\mathbf{y}) = g(\mathbf{y}) - p_\ell f_\ell(\mathbf{y}), \quad \mathbf{y} \in \mathbb{R}^k$$

and the test $d^* : \mathbb{R}^k \rightarrow \{0, 1, \dots, M-1\}$ can be rewritten more compactly as

$$d^*(\mathbf{y}) = \arg \max (\ell = 0, \dots, M-1 : p_\ell f_\ell(\mathbf{y})), \quad \mathbf{y} \in \mathbb{R}^k \quad (7.20)$$

with a lexicographic tiebreaker in the event of ties.

The ML test When all the hypotheses are *equally likely*, namely

$$p_0 = \dots = p_{M-1} = \frac{1}{M},$$

then (7.20) becomes

$$d^*(\mathbf{y}) = \arg \max (\ell = 0, \dots, M-1 : f_\ell(\mathbf{y})), \quad \mathbf{y} \in \mathbb{R}^k \quad (7.21)$$

with a lexicographic tiebreaker in the event of ties, so that

$$d^*(\mathbf{y}) = m \quad \text{iff} \quad f_m(\mathbf{y}) = \max (f_\ell(\mathbf{y}), \ell = 0, 1, \dots, M-1) \quad (7.22)$$

with a lexicographic tiebreaker in the event of ties.

The MAP computer Bayes's Theorem gives

$$\mathbb{P}[H = \ell | \mathbf{Y} = \mathbf{y}] = \frac{p_\ell f_\ell(\mathbf{y})}{\sum_{m=0}^{M-1} p_m f_m(\mathbf{y})}, \quad \ell = 0, 1, \dots, M-1, \quad \mathbf{y} \in \mathbb{R}^k.$$

This allows rewriting the test $d^* : \mathbb{R}^k \rightarrow \{0, 1, \dots, M-1\}$ in the more compact form

$$d^*(\mathbf{y}) = \arg \max (\ell = 0, \dots, M-1 : \mathbb{P}[H = \ell | \mathbf{Y} = \mathbf{y}]), \quad \mathbf{y} \in \mathbb{R}^k \quad (7.23)$$

with a lexicographic tiebreaker in the event of ties. As with the binary case, we refer to this rule as the *Maximum A Posteriori* computer.

7.7 The Gaussian case

Consider the case where the rv \mathbf{Y} is Gaussian under each hypothesis, namely

$$H_m : \mathbf{Y} \sim N(\mathbf{a}_m, \mathbf{R}_m). \quad m = 0, 1, \dots, M-1. \quad (7.24)$$

where $\mathbf{a}_0, \dots, \mathbf{a}_{M-1}$ are elements in \mathbb{R}^k , and the $k \times k$ symmetric matrices $\mathbf{R}_0, \dots, \mathbf{R}_{M-1}$ are *positive definite* (thus *invertible*). Condition **(A.3)** holds with respect to Lebesgue measure.

Throughout the M pairs $(\mathbf{a}_0, \mathbf{R}_0), \dots, (\mathbf{a}_{M-1}, \mathbf{R}_{M-1})$ are assumed to be distinct so that the probability density function $f_0, \dots, f_{M-1} : \mathbb{R}^k \rightarrow \mathbb{R}_+$ are distinct. Indeed, for each $m = 0, \dots, M-1$, the probability density function $f_m : \mathbb{R}^k \rightarrow \mathbb{R}_+$ is given by

$$f_m(\mathbf{y}) = \frac{1}{\sqrt{(2\pi)^k \det \mathbf{R}_m}} e^{-\frac{1}{2} Q_m(\mathbf{y})}, \quad \mathbf{y} \in \mathbb{R}^k$$

with

$$Q_m(\mathbf{y}) = (\mathbf{y} - \mathbf{a}_m)' \mathbf{R}_m^{-1} (\mathbf{y} - \mathbf{a}_m), \quad \mathbf{y} \in \mathbb{R}^k.$$

Thus, for each $\ell = 0, \dots, M-1$, we have

$$\begin{aligned} & h_\ell(\mathbf{y}) \\ &= \sum_{m=0, m \neq \ell}^{M-1} p_m (C(m, \ell) - C(m, m)) \frac{1}{\sqrt{(2\pi)^k \det \mathbf{R}_m}} e^{-\frac{1}{2} Q_m(\mathbf{y})}, \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned}$$

For the probability of error criterion, for each $\ell = 0, \dots, M - 1$, this last expression becomes

$$\begin{aligned} h_\ell(\mathbf{y}) &= \sum_{m=0, m \neq \ell}^{M-1} \frac{p_m}{\sqrt{(2\pi)^k \det \mathbf{R}_m}} e^{-\frac{1}{2}Q_m(\mathbf{y})} \\ &= g(\mathbf{y}) - \frac{p_\ell}{\sqrt{(2\pi)^k \det \mathbf{R}_\ell}} e^{-\frac{1}{2}Q_\ell(\mathbf{y})}, \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned} \quad (7.25)$$

The equal covariance case When

$$\mathbf{R}_0 = \dots = \mathbf{R}_{M-1} \equiv \mathbf{R},$$

additional simplifications occur: For each $\ell = 0, \dots, M - 1$, we find

$$\begin{aligned} h_\ell(\mathbf{y}) &= \frac{1}{\sqrt{(2\pi)^k \det \mathbf{R}}} \sum_{m=0, m \neq \ell}^{M-1} \frac{p_m (C(m, \ell) - C(m, m))}{\sqrt{(2\pi)^k \det \mathbf{R}}} e^{-\frac{1}{2}Q(\mathbf{y} - \mathbf{a}_m)}, \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned}$$

with

$$Q(\mathbf{y}) = \mathbf{y}' \mathbf{R}^{-1} \mathbf{y}, \quad \mathbf{y} \in \mathbb{R}^k.$$

For the probability of error criterion, for each $\ell = 0, \dots, M - 1$, we find

$$\begin{aligned} h_\ell(\mathbf{y}) &= \sum_{m=0, m \neq \ell}^{M-1} \frac{p_m}{\sqrt{(2\pi)^k \det \mathbf{R}}} e^{-\frac{1}{2}Q(\mathbf{y} - \mathbf{a}_m)} \\ &= g(\mathbf{y}) - \frac{p_\ell}{\sqrt{(2\pi)^k \det \mathbf{R}}} e^{-\frac{1}{2}Q(\mathbf{y} - \mathbf{a}_\ell)}, \quad \mathbf{y} \in \mathbb{R}^k. \end{aligned} \quad (7.26)$$

Writing

$$d(p; \mathbf{y}) = \log p - \frac{1}{2}Q(\mathbf{y}), \quad \begin{array}{l} p \in (0, 1) \\ \mathbf{y} \in \mathbb{R}^k, \end{array}$$

the Bayesian test $d^* : \mathbb{R}^k \rightarrow \{0, \dots, M - 1\}$ now reduces to

$$d^*(\mathbf{y}) = m \quad \text{iff} \quad \begin{array}{l} d(p_m; \mathbf{y} - \mathbf{a}_m) \\ = \max (d(p_\ell; \mathbf{y} - \mathbf{a}_\ell), \ell = 0, 1, \dots, M - 1) \end{array}$$

with a lexicographic tiebreaker in the event of ties.

7.8 Exercises

Chapter 8

Composite hypothesis testing problems

8.1 A motivating example

Consider the following problem of deciding between two hypotheses. Under the null hypothesis H_0 , the observation \mathbf{Y} is an \mathbb{R}^k -valued rv which is normally distributed with mean vector \mathbf{m} and covariance matrix \mathbf{R} which are both *known*. Under the alternative hypothesis H_1 , the \mathbb{R}^k -valued rv \mathbf{Y} is normally distributed with mean vector $\boldsymbol{\theta}$ and covariance matrix \mathbf{R} where $\boldsymbol{\theta} \neq \mathbf{m}$ is only known to lie in a subset Θ_1 of \mathbb{R}^k , and is otherwise unspecified. We assume that \mathbf{m} is *not* an element of Θ_1 .

This problem of testing for the binary hypothesis H_0 vs. H_1 can also be interpreted as one of deciding between the hypotheses

$$\begin{aligned} H_1 &: \{H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1\} \\ H_0 &: \mathbf{Y} \sim N(\mathbf{m}, \mathbf{R}). \end{aligned} \tag{8.1}$$

where for each $\boldsymbol{\theta} \in \mathbb{R}^k$, we write

$$H_{\boldsymbol{\theta}} : \mathbf{Y} \sim N(\boldsymbol{\theta}, \mathbf{R}).$$

In such situations, when Θ_1 is not reduced to a singleton, the alternative hypothesis can be viewed as a *composite* hypothesis $\{H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1\}$. We emphasize that we seek to decide between H_0 and H_1 , or equivalently, between H_0 and $\{H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1\}$; the precise value of $\boldsymbol{\theta}$ is not thought

8.2 The probabilistic model

More generally consider two non-empty Borel subsets Θ_0 and Θ_1 of \mathbb{R}^p for some positive integer p . Assume that

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

We shall set

$$\Theta = \Theta_0 \cup \Theta_1.$$

so that the pair Θ_0 and Θ_1 forms a Borel partition of Θ .

Given is a family of probability distributions $\{F_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta\}$ on \mathbb{R}^k . For mathematical reasons, it is required that the mapping $\Theta \times \mathbb{R}^k \rightarrow [0, 1] : (\boldsymbol{\theta}, \mathbf{y}) \rightarrow F_{\boldsymbol{\theta}}(\mathbf{y})$ be Borel measurable. This condition is satisfied in all applications of interest.

We are given a measurable space (Ω, \mathcal{F}) . The state of nature is modeled by means of a rv $\vartheta : \Omega \rightarrow \Theta$ defined on (Ω, \mathcal{F}) . The observation is given by an \mathbb{R}^k -valued rv $\mathbf{Y} : \Omega \rightarrow \mathbb{R}^k$ defined on the same measurable space (Ω, \mathcal{F}) . with the interpretation that

$$\mathbb{P}[\mathbf{Y} \leq \mathbf{y} | \vartheta = \boldsymbol{\theta}] = F_{\boldsymbol{\theta}}(\mathbf{y}), \quad \begin{array}{l} \mathbf{y} \in \mathbb{R}^k, \\ \boldsymbol{\theta} \in \Theta. \end{array} \quad (8.2)$$

The state of nature and the corresponding observation are summarized as

$$H_{\boldsymbol{\theta}} : \mathbf{Y} \sim F_{\boldsymbol{\theta}} \quad (8.3)$$

with $\boldsymbol{\theta}$ ranging in Θ .

The *composite* binary hypothesis testing problem is then the problem of deciding between the two composite hypotheses $H_0 = \{H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_0\}$ and $H_1 = \{H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1\}$ on the basis of the observation \mathbf{Y} .

If either Θ_0 or Θ_1 is reduced to a single element, the corresponding hypothesis is termed *simple*. Obviously the problems of simple binary hypothesis testing discussed in Chapters 1–5 obtains when each of the sets Θ_0 and Θ_1 contains exactly one element.

The composite binary hypothesis testing problem is concisely denoted by

$$\begin{array}{l} H_1 : \mathbf{Y} \sim F_{\boldsymbol{\theta}}, \quad \boldsymbol{\theta} \in \Theta_1 \\ H_0 : \mathbf{Y} \sim F_{\boldsymbol{\theta}}, \quad \boldsymbol{\theta} \in \Theta_0 \end{array} \quad (8.4)$$

Sometimes, once the family of probability distributions $\{F_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta\}$ has been specified, the notation is simplified to read

$$\begin{array}{l} H_1 : \boldsymbol{\theta} \in \Theta_1 \\ H_0 : \boldsymbol{\theta} \in \Theta_0 \end{array} \quad (8.5)$$

As in earlier chapters, the discussion will require that certain assumptions are enforced. Condition (A.4) given next parallels Condition (A.1) given in the binary case and Condition (A.3) given in the M -ary case.

Condition (A.4): For each θ in Θ , the probability distribution F_θ on \mathbb{R}^k is absolutely continuous with respect to some distribution F on \mathbb{R}^k .

This is equivalent to requiring that for each θ in Θ , there exists a Borel mapping $f_\theta : \mathbb{R}^k \rightarrow \mathbb{R}_+$ such that

$$F_\theta(\mathbf{y}) = \int_{-\infty}^{\mathbf{y}} f_\theta(\boldsymbol{\eta}) dF(\boldsymbol{\eta}), \quad \mathbf{y} \in \mathbb{R}^k. \quad (8.6)$$

For mathematical reasons we require that the mapping

$$\Theta \times \mathbb{R}^k \rightarrow \mathbb{R}_+ : (\theta, \mathbf{y}) \rightarrow f_\theta(\mathbf{y})$$

be Borel measurable. This condition is satisfied in all applications of interest.

In a manner reminiscent of the parameter estimation problem of Chapter 1 there are two possible cases, depending on whether or not θ is modeled as a rv; these are the Bayesian and non-Bayesian cases, respectively.

8.3 The Bayesian case

In some settings it is appropriate to imagine that the value of the parameter θ is indeed through some randomization mechanism. Thus, assume that there exists a Θ -valued rv ϑ defined on the measurable space (Ω, \mathcal{F}) , and let $K : \mathbb{R}^p \rightarrow [0, 1]$ denote its probability distribution. Thus,

$$\mathbb{P}[\vartheta \leq \mathbf{t}] = K(\mathbf{t}), \quad \mathbf{t} \in \mathbb{R}^p.$$

The requirement that $\mathbb{P}[\vartheta \in \Theta] = 1$ is equivalent to

$$\int_{\Theta^c} dK(\mathbf{t}) = 0.$$

As we now show, this composite binary hypothesis testing problem can be reformulated as a simple binary hypothesis testing problem. To do so we define the $\{0, 1\}$ -valued rv H given by

$$H \equiv \mathbf{1}[\vartheta \in \Theta_1].$$

Note that

$$\begin{aligned}
 p &= \mathbb{P}[H = 1] \\
 &= \mathbb{P}[\vartheta \in \Theta_1] \\
 &= \int_{\Theta_1} dK(\mathbf{t}).
 \end{aligned} \tag{8.7}$$

In a similar manner, since $1 - H = \mathbf{1}[\vartheta \in \Theta_0]$, we conclude that

$$\begin{aligned}
 1 - p &= \mathbb{P}[H = 0] \\
 &= \mathbb{P}[\vartheta \in \Theta_0] \\
 &= \int_{\Theta_0} dK(\mathbf{t}).
 \end{aligned} \tag{8.8}$$

For each $h = 0, 1$, the conditional probability distribution of \mathbf{Y} given that $H = h$ can be calculated as

$$\begin{aligned}
 F_h(\mathbf{y}) &\equiv \mathbb{P}[\mathbf{Y} \leq \mathbf{y} | H = h] \\
 &= \frac{\mathbb{P}[\mathbf{Y} \leq \mathbf{y}, H = h]}{\mathbb{P}[H = h]} \\
 &= \frac{\mathbb{P}[\mathbf{Y} \leq \mathbf{y}, \vartheta \in \Theta_h]}{\mathbb{P}[\vartheta \in \Theta_h]} \\
 &= \frac{\mathbb{E}[\mathbb{E}[\mathbf{1}[\mathbf{Y} \leq \mathbf{y}] | \vartheta] \mathbf{1}[\vartheta \in \Theta_h]]}{\mathbb{P}[\vartheta \in \Theta_h]} \\
 &= \frac{\mathbb{E}[\mathbb{P}[\mathbf{Y} \leq \mathbf{y} | \vartheta] \mathbf{1}[\vartheta \in \Theta_h]]}{\mathbb{P}[\vartheta \in \Theta_h]} \\
 &= \frac{\mathbb{E}[F_\vartheta(\mathbf{y}) \mathbf{1}[\vartheta \in \Theta_h]]}{\mathbb{P}[\vartheta \in \Theta_h]} \\
 &= \frac{\int_{\Theta_h} F_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_h} dK(\mathbf{t})}, \quad \mathbf{y} \in \mathbb{R}^k.
 \end{aligned} \tag{8.9}$$

Fix $h = 0, 1$ and \mathbf{y} in \mathbb{R}^k . Under Assumption **(A.4)** we note that

$$\begin{aligned}
 \int_{\Theta_h} F_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) &= \int_{\Theta_h} \left(\int_{-\infty}^{\mathbf{y}} f_{\mathbf{t}}(\boldsymbol{\eta}) dF(\boldsymbol{\eta}) \right) dK(\mathbf{t}) \\
 &= \int_{-\infty}^{\mathbf{y}} \left(\int_{\Theta_h} f_{\mathbf{t}}(\boldsymbol{\eta}) dK(\mathbf{t}) \right) dF(\boldsymbol{\eta})
 \end{aligned} \tag{8.10}$$

by Tonelli's Theorem, whence

$$F_h(\mathbf{y}) = \int_{-\infty}^{\mathbf{y}} \left(\frac{\int_{\Theta_h} f_{\mathbf{t}}(\boldsymbol{\eta}) dK(\mathbf{t})}{\int_{\Theta_h} dK(\mathbf{t})} \right) dF(\boldsymbol{\eta}). \quad (8.11)$$

Thus, Condition **(A.1)** holds for F_0 and F_1 with respect to F with probability density functions $f_0, f_1 : \mathbb{R}^k \rightarrow \mathbb{R}_+$ given by

$$f_h(\mathbf{y}) \equiv \frac{\int_{\Theta_h} f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_h} dK(\mathbf{t})}, \quad h = 0, 1, \quad \mathbf{y} \in \mathbb{R}^k.$$

At this point the reader might expect the corresponding tests $\{d_\eta, \eta \geq 0\}$ to play a prominent role where for each $\eta \geq 0$, the test $d_\eta : \mathbb{R}^k \rightarrow \{0, 1\}$ is given by

$$\begin{aligned} d_\eta(\mathbf{y}) = 0 & \quad \text{iff} \quad f_1(\mathbf{y}) < \eta f_0(\mathbf{y}) \\ & \quad \text{iff} \quad \frac{\int_{\Theta_1} f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_1} dK(\mathbf{t})} < \eta \frac{\int_{\Theta_0} f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_0} dK(\mathbf{t})}. \end{aligned} \quad (8.12)$$

In Section 8.9 we will see that this is not the case.

8.4 The Bayesian cost problem

In the Bayesian setup described in Section 8.3, we proceed as in Chapter 2 by introducing a cost incurred for making decisions. This is quantified by the mapping $C : \Theta \times \{0, 1\} \rightarrow \mathbb{R}$ with the interpretation that

$$C(\boldsymbol{\theta}, d) = \begin{array}{l} \text{Cost incurred for deciding } d \\ \text{when } \boldsymbol{\vartheta} = \boldsymbol{\theta} \end{array}, \quad \begin{array}{l} \boldsymbol{\theta} \in \Theta \\ d = 0, 1. \end{array}$$

We require that for each $d = 0, 1$, the mapping $\Theta \rightarrow \mathbb{R} : \boldsymbol{\theta} \rightarrow C(\boldsymbol{\theta}, d)$ is Borel measurable. This guarantees that for every test $d : \mathbb{R}^k \rightarrow \{0, 1\}$, $C(\boldsymbol{\vartheta}, d(\mathbf{Y}))$ is a rv defined on (Ω, \mathcal{F}) . To avoid unnecessary technical difficulties (and for ease of exposition) we further assume that

$$0 \leq C(\boldsymbol{\theta}, d) \leq B, \quad \begin{array}{l} d = 0, 1 \\ \boldsymbol{\vartheta} \in \Theta \end{array}$$

for some scalar B . Together these requirements ensure that the expectation $\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))]$ is well defined and finite.

With any admissible test $d : \mathbb{R}^k \rightarrow \{0, 1\}$, we define the expected cost

$$J(d) = \mathbb{E} [C(\boldsymbol{\vartheta}, d(\mathbf{Y}))]. \quad (8.13)$$

As before, the *Bayesian Problem* \mathcal{P}_B is the minimization problem

$$\mathcal{P}_B : \quad \text{Minimize } J(d) \text{ over } d \text{ in } \mathcal{D}.$$

This amounts to finding an admissible test d^* in \mathcal{D} such that

$$J(d^*) \leq J(d), \quad d \in \mathcal{D}. \quad (8.14)$$

Any admissible test d^* which satisfies (8.14) is called a Bayesian test, and the value

$$J(d^*) = \inf_{d \in \mathcal{D}} J(d) = \min_{d \in \mathcal{D}} J(d) \quad (8.15)$$

is known as the *Bayesian cost*.

In Section 8.3 we have shown that the Bayesian formulation of the composite binary hypothesis problem can be recast as a simple binary hypothesis problem. However, as discussed in Section 8.9, in general it is not possible to write

$$J(d) = \mathbb{E} [C_{\text{New}}(H, d(\mathbf{Y}))], \quad d \in \mathcal{D}$$

for some mapping $C_{\text{New}} : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}$ (possibly derived from the original cost function $C : \Theta \times \{0, 1\} \rightarrow \mathbb{R}$). This already suggests that the Bayesian test may not belong to the class of tests $\{d_\eta, \eta \geq 0\}$ introduced at (8.12).

Solving the Bayesian Problem \mathcal{P}_B In view of these comments a different approach is needed. Fix d in \mathcal{D} . A standard preconditioning argument gives

$$\begin{aligned} J(d) &= \mathbb{E} [C(\boldsymbol{\vartheta}, d(\mathbf{Y}))] \\ &= \mathbb{E} [\mathbb{E} [C(\boldsymbol{\vartheta}, d(\mathbf{Y})) | \mathbf{Y}]] \\ &= \mathbb{E} [\mathbb{E} [\mathbf{1} [d(\mathbf{Y}) = 0] C(\boldsymbol{\vartheta}, d(\mathbf{Y})) + \mathbf{1} [d(\mathbf{Y}) = 1] C(\boldsymbol{\vartheta}, d(\mathbf{Y})) | \mathbf{Y}]] \\ &= \mathbb{E} [\mathbb{E} [\mathbf{1} [d(\mathbf{Y}) = 0] C(\boldsymbol{\vartheta}, 0) + \mathbf{1} [d(\mathbf{Y}) = 1] C(\boldsymbol{\vartheta}, 1) | \mathbf{Y}]] \\ &= \mathbb{E} [\mathbf{1} [d(\mathbf{Y}) = 0] \mathbb{E} [C(\boldsymbol{\vartheta}, 0) | \mathbf{Y}] + \mathbf{1} [d(\mathbf{Y}) = 1] \mathbb{E} [C(\boldsymbol{\vartheta}, 1) | \mathbf{Y}]] \\ &= \mathbb{E} [\mathbf{1} [d(\mathbf{Y}) = 0] \widehat{C}(0, \mathbf{Y}) + \mathbf{1} [d(\mathbf{Y}) = 1] \widehat{C}(1, \mathbf{Y})] \end{aligned} \quad (8.16)$$

where we have introduced the Borel mapping $\widehat{C} : \{0, 1\} \times \mathbb{R}^k \rightarrow \mathbb{R}$ given by

$$\widehat{C}(d, \mathbf{y}) \equiv \mathbb{E}[C(\boldsymbol{\vartheta}, d) | \mathbf{Y} = \mathbf{y}], \quad \begin{array}{l} d = 0, 1 \\ \mathbf{y} \in \mathbb{R}^k. \end{array}$$

Theorem 8.4.1 *Under the foregoing assumptions, the test $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ defined by*

$$d^*(\mathbf{y}) = \begin{cases} 0 & \text{if } \widehat{C}(0, \mathbf{y}) < \widehat{C}(1, \mathbf{y}) \\ 1 & \text{if } \widehat{C}(1, \mathbf{y}) \leq \widehat{C}(0, \mathbf{y}) \end{cases} \quad (8.17)$$

is admissible and solves the Bayesian Problem \mathcal{P}_B .

Note that the result does not even depend on Assumption (A.4). Moreover, when $\widehat{C}(0, \mathbf{y}) = \widehat{C}(1, \mathbf{y})$, we may select $d^*(\mathbf{y}) = 0$ or $d^*(\mathbf{y}) = 1$ somewhat arbitrarily as long as the resulting mapping $\mathbb{R}^k \rightarrow \{0, 1\}$ is Borel measurable.

Proof. The admissibility of d^* follows from the Borel measurability of the mapping $\widehat{C} : \{0, 1\} \times \mathbb{R}^k \rightarrow \mathbb{R}$. To show its optimality note that for each test d in \mathcal{D} it holds that

$$\begin{aligned} J(d) &= \mathbb{E} \left[\mathbf{1}[d(\mathbf{Y}) = 0] \widehat{C}(0, \mathbf{Y}) + \mathbf{1}[d(\mathbf{Y}) = 1] \widehat{C}(1, \mathbf{Y}) \right] \\ &= \mathbb{E} \left[\widehat{C}(1, \mathbf{Y}) \right] + \mathbb{E} \left[\mathbf{1}[d(\mathbf{Y}) = 0] \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \\ &= \mathbb{E} \left[\widehat{C}(1, \mathbf{Y}) \right] + \mathbb{E} \left[\mathbf{1}[\mathbf{Y} \in C(d)] \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \end{aligned}$$

where as usual the Borel set $C(d)$ is given by (1.12). The conclusion readily follows by arguments similar to the ones used for the proof of Theorem 2.2.1: Indeed, we have

$$\begin{aligned} J(d) - J(d^*) &= \mathbb{E} \left[(\mathbf{1}[\mathbf{Y} \in C(d)] - \mathbf{1}[\mathbf{Y} \in C(d^*)]) \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \\ &= \mathbb{E} \left[\mathbf{1}[\mathbf{Y} \in C(d) \setminus C(d^*)] \cdot \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \\ &\quad - \mathbb{E} \left[\mathbf{1}[\mathbf{Y} \in C(d^*) \setminus C(d)] \cdot \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \end{aligned}$$

As we note that

$$C(d^*) \equiv \left\{ \mathbf{y} \in \mathbb{R}^k : \widehat{C}(0, \mathbf{y}) < \widehat{C}(1, \mathbf{y}) \right\}$$

it is plain that

$$\mathbb{E} \left[\mathbf{1} [\mathbf{Y} \in C(d) \setminus C(d^*)] \cdot \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \geq 0$$

and

$$\mathbb{E} \left[\mathbf{1} [\mathbf{Y} \in C(d^*) \setminus C(d)] \cdot \left(\widehat{C}(0, \mathbf{Y}) - \widehat{C}(1, \mathbf{Y}) \right) \right] \leq 0.$$

It follows that $J(d^*) \leq J(d)$ and the optimality of d^* is now established. \blacksquare

8.5 The non-Bayesian case – A generalized Neyman-Pearson formulation

In this formulation θ is an unknown parameter lying in Θ , and an approach à la Neyman-Pearson seems warranted. Since composite hypotheses are now involved, earlier definitions given in Chapter 4 need to be modified.

Consider a test $d : \mathbb{R}^k \rightarrow \{0, 1\}$ in \mathcal{D} that tests the null hypothesis H_0 against the alternative H_1 . We define its *size* to be the quantity

$$\alpha_{\Theta_0}(d) \equiv \sup_{\boldsymbol{\theta} \in \Theta_0} \mathbb{P}_{\boldsymbol{\theta}} [d(\mathbf{Y}) = 1]. \quad (8.18)$$

With $\boldsymbol{\theta}$ in Θ_0 , the probability $\mathbb{P}_{\boldsymbol{\theta}} [d(\mathbf{Y}) = 0]$ can be interpreted as the probability of false alarm under the test d given that the hypothesis $H_{\boldsymbol{\theta}}$ is indeed correct.

Fix α in $[0, 1]$. Let $\mathcal{D}_{\Theta_0, \alpha}$ denote the collection of all tests in \mathcal{D} whose size is no greater than α , namely,

$$\mathcal{D}_{\Theta_0, \alpha} \equiv \{d \in \mathcal{D} : \alpha_{\Theta_0}(d) \leq \alpha\}. \quad (8.19)$$

When Θ_0 is reduced to a singleton $\boldsymbol{\theta}_0$ we write $\mathcal{D}_{\boldsymbol{\theta}_0, \alpha}$ instead of $\mathcal{D}_{\{\boldsymbol{\theta}_0\}, \alpha}$, so that

$$\mathcal{D}_{\boldsymbol{\theta}_0, \alpha} \equiv \{d \in \mathcal{D} : \mathbb{P}_{\boldsymbol{\theta}_0} [d(\mathbf{Y}) = 1] \leq \alpha\}.$$

The inclusion

$$\mathcal{D}_{\Theta_0, \alpha} \subseteq \mathcal{D}_{\theta_0, \alpha}, \quad \theta_0 \in \Theta_0 \quad (8.20)$$

always holds – It is also easy to see that

$$\mathcal{D}_{\Theta_0, \alpha} = \bigcap_{\theta_0 \in \Theta_0} \mathcal{D}_{\theta_0, \alpha}.$$

When dealing with composite hypotheses, the generalized Neyman-Pearson formulation takes the following form:

For α in $[0, 1]$, find a test $d_{\text{UMP}}(\alpha) : \mathbb{R}^k \rightarrow \{0, 1\}$ in $\mathcal{D}_{\Theta_0, \alpha}$ which is optimal in the sense that

$$\mathbb{P}_{\theta} [d(\mathbf{Y}) = 1] \leq \mathbb{P}_{\theta} [d_{\text{UMP}}(\alpha)(\mathbf{Y}) = 1], \quad \begin{array}{l} \theta \in \Theta_1, \\ d \in \mathcal{D}_{\Theta_0, \alpha}. \end{array} \quad (8.21)$$

Such a test $d_{\text{UMP}}(\alpha)$, when it exists, is called a *Uniformly Most Powerful (UMP)* test of size α for testing for testing $H_0 : H_{\theta}$, $\theta \in \Theta_0$ against $H_1 : H_{\theta}$, $\theta \in \Theta_1$,

A natural question is whether UMP tests exist and when they do, how does one go about identifying them. The discussion in Section 8.6 below provides some useful pointers concerning these issues.

8.6 Searching for UMP tests

If a UMP test $d_{\text{UMP}}(\alpha)$ of size α did exist for testing $H_0 : H_{\theta}$, $\theta \in \Theta_0$ against $H_1 : H_{\theta}$, $\theta \in \Theta_1$, then by definition $d_{\text{UMP}}(\alpha)$ is a test in $\mathcal{D}_{\Theta_0, \alpha}$, hence in $\mathcal{D}_{\theta_0, \alpha}$ for each θ_0 in Θ_0 . However, in general the optimality property (8.21) does not necessarily imply

$$\mathbb{P}_{\theta} [d(\mathbf{Y}) = 1] \leq \mathbb{P}_{\theta} [d_{\text{UMP}}(\alpha)(\mathbf{Y}) = 1], \quad \begin{array}{l} \theta \in \Theta_1, \\ d \in \mathcal{D}_{\theta_0, \alpha} \end{array} \quad (8.22)$$

since *a priori* the inclusion (8.20) may be strict.

Were it the case that $\mathcal{D}_{\Theta_0, \alpha} = \mathcal{D}_{\theta_0, \alpha}$ for some θ_0 in Θ_0 , then (8.22) would hold (as it now coincides with (8.21)), and the test $d_{\text{UMP}}(\alpha)$ would therefore act as a Neyman-Pearson test $d_{\text{NP}}(\alpha; \theta_0, \theta)$ of size α for testing $H_0 \equiv H_{\theta_0}$ against $H_1 \equiv H_{\theta}$ for each θ in Θ_1 .

There is of course one situation when the equality $\mathcal{D}_{\Theta_0, \alpha} = \mathcal{D}_{\theta_0, \alpha}$ obviously holds, namely when H_0 is a simple hypothesis so that Θ_0 is reduced to a singleton,

say $\Theta_0 = \{\boldsymbol{\theta}_0\}$ for some $\boldsymbol{\theta}_0$ not in Θ_1 . The discussion above then implies that any UMP test $d_{\text{UMP}}(\alpha)$, if one exists, for testing $H_0 \equiv H_{\boldsymbol{\theta}_0}$ against $H_1 : H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1$, must satisfy

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\theta}} [d_{\text{UMP}}(\alpha)(\mathbf{Y}) = 1] \\ &= \mathbb{P}_{\boldsymbol{\theta}} [d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})(\mathbf{Y}) = 1], \quad \boldsymbol{\theta} \in \Theta_1. \end{aligned} \quad (8.23)$$

This is clearly a non-trivial restriction on the problem, and already suggests that UMP tests may not always exist even H_0 is a simple hypothesis – This is discussed on an example in Section 8.8.

Nevertheless, when H_0 is a simple hypothesis, these observations do point to an obvious strategy for finding UMPs: For each $\boldsymbol{\theta}$ in Θ_1 , determine the Neyman-Pearson test $d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})$ of size α for testing $H_0 \equiv H_{\boldsymbol{\theta}_0}$ against $H_1 \equiv H_{\boldsymbol{\theta}}$. If its implementation does not require explicit knowledge of $\boldsymbol{\theta}$, then the set $C(d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta}))$ will be independent of $\boldsymbol{\theta}$ in the sense that there exists a Borel subset C of \mathbb{R}^k such that $C(d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})) = C$ for every $\boldsymbol{\theta}$ in Θ_1 . The test $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ defined by

$$d^*(\mathbf{y}) = \begin{cases} 0 & \text{if } \mathbf{y} \in C \\ 1 & \text{if } \mathbf{y} \notin C \end{cases}$$

is an admissible test in $\mathcal{D}_{\boldsymbol{\theta}_0, \alpha}$ since by construction

$$\begin{aligned} \mathbb{P}_{\boldsymbol{\theta}_0} [d^*(\mathbf{Y}) = 1] &= \mathbb{P}_{\boldsymbol{\theta}_0} [\mathbf{Y} \notin C] \\ &= \mathbb{P}_{\boldsymbol{\theta}_0} [\mathbf{Y} \notin C(d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta}))] \\ &= \mathbb{P}_{\boldsymbol{\theta}_0} [d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})(\mathbf{Y}) = 1] = \alpha. \end{aligned} \quad (8.24)$$

By the same arguments we also conclude that

$$\mathbb{P}_{\boldsymbol{\theta}} [d^*(\mathbf{Y}) = 1] = \mathbb{P}_{\boldsymbol{\theta}} [d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})(\mathbf{Y}) = 1], \quad \boldsymbol{\theta} \in \Theta_1 \quad (8.25)$$

and we conclude that d^* is a UMP test for of size α for testing for testing $H_0 \equiv H_{\boldsymbol{\theta}_0}$ against $H_1 : H_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta_1$,

When H_0 is a composite hypothesis, finding a UMP test can be quite tricky. A first natural step would consist in finding the Neyman-Pearson test $d_{\text{NP}}(\alpha; \boldsymbol{\theta}_0, \boldsymbol{\theta})$ of size α to test $H_0 \equiv H_{\boldsymbol{\theta}_0}$ against $H_1 \equiv H_{\boldsymbol{\theta}}$ with $\boldsymbol{\theta}_0$ and $\boldsymbol{\theta}_1$ arbitrary in Θ_0 and Θ_1 , respectively. In some cases exploring the structure of these tests may lead to the UMP test of size α .

These ideas are illustrated through an example in Section 8.7 and Section 8.8.

8.7 An example

The discussion of Section 8.5 will be illustrated in the case when $\Theta \subseteq \mathbb{R}$ and the probability distributions $\{F_\theta, \theta \in \mathbb{R}\}$ are unit-variance Gaussian distributions on \mathbb{R} : Thus, with θ arbitrary in \mathbb{R} ,

$$H_\theta : Y \sim N(\theta, 1)$$

so that F_θ admits the density $f_\theta : \mathbb{R} \rightarrow \mathbb{R}_+$ given by

$$f_\theta(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\theta)^2}, \quad y \in \mathbb{R}.$$

With distinct θ_0 and θ_1 in \mathbb{R} , consider the Neyman-Pearson formulation for the binary hypothesis problem

$$\begin{aligned} H_1 : Y &\sim N(\theta_1, 1) \\ H_0 : Y &\sim N(\theta_0, 1) \end{aligned} \quad (8.26)$$

As shown in Section 5.5 this problem has a complete solution: Fix $\lambda > 0$. The test $d_\lambda : \mathbb{R} \rightarrow \{0, 1\}$ takes the form

$$\begin{aligned} d_\lambda(y) = 0 & \quad \text{iff} \quad \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\theta_1)^2} < \lambda \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\theta_0)^2} \\ & \quad \text{iff} \quad (y - \theta_0)^2 < 2 \log \lambda + (y - \theta_1)^2 \\ & \quad \text{iff} \quad 2y(\theta_1 - \theta_0) < 2 \log \lambda + \theta_1^2 - \theta_0^2. \end{aligned} \quad (8.27)$$

For notational convenience we shall write

$$T_\lambda(\theta_0; \theta_1) \equiv \frac{\log \lambda}{\theta_1 - \theta_0} + \frac{\theta_1 + \theta_0}{2}.$$

Two cases arise: If $\theta_0 < \theta_1$, then

$$d_\lambda(y) = 0 \quad \text{iff} \quad y < T_\lambda(\theta_0; \theta_1) \quad (8.28)$$

and by standard arguments we get

$$\begin{aligned} \mathbb{P}_{\theta_0} [d_\lambda(\mathbf{Y}) = 1] &= \mathbb{P}_{\theta_0} [Y \geq T_\lambda(\theta_0; \theta_1)] \\ &= \mathbb{P}_{\theta_0} \left[Y - \theta_0 \geq \frac{\log \lambda}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} \right] \\ &= 1 - \Phi \left(\frac{\log \lambda}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} \right). \end{aligned} \quad (8.29)$$

If $\theta_1 < \theta_0$, then

$$d_\lambda(y) = 0 \quad \text{iff} \quad y > T_\lambda(\theta_0; \theta_1), \quad (8.30)$$

and this time we find

$$\begin{aligned} \mathbb{P}_{\theta_0} [d_\lambda(\mathbf{Y}) = 1] &= \mathbb{P}_{\theta_0} [Y \leq T_\lambda(\theta_0; \theta_1)] \\ &= \mathbb{P}_{\theta_0} \left[Y - \theta_0 \leq \frac{\log \lambda}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} \right] \\ &= \Phi \left(\frac{\log \lambda}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} \right). \end{aligned} \quad (8.31)$$

Fix α in $(0, 1)$. The Neyman-Pearson test of size α for testing H_{θ_1} against H_{θ_0} is the test $d_{\lambda(\theta_1, \theta_0; \alpha)}$ where $\lambda(\theta_1, \theta_0; \alpha)$ is that value of $\lambda > 0$ determined by $\mathbb{P}_{\theta_0} [d_\lambda(\mathbf{Y}) = 1] = \alpha$.

If $\theta_0 < \theta_1$, then

$$\frac{\log \lambda(\theta_1, \theta_0; \alpha)}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} = \Phi^{-1}(1 - \alpha),$$

and the test $d_{\lambda(\theta_1, \theta_0; \alpha)}$ is given by

$$d_{\lambda(\theta_1, \theta_0; \alpha)}(y) = 0 \quad \text{iff} \quad y < \theta_0 + \Phi^{-1}(1 - \alpha).$$

If $\theta_1 < \theta_0$, then

$$\frac{\log \lambda(\theta_1, \theta_0; \alpha)}{\theta_1 - \theta_0} + \frac{\theta_1 - \theta_0}{2} = \Phi^{-1}(\alpha),$$

and the test $d_{\lambda(\theta_1, \theta_0; \alpha)}$ is given by

$$d_{\lambda(\theta_1, \theta_0; \alpha)}(y) = 0 \quad \text{iff} \quad y > \theta_0 + \Phi^{-1}(\alpha).$$

8.8 UMP tests for the example

We now consider four different situations, each associated with different sets Θ_0 and Θ_1 .

Case I: $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = (\theta_0, \infty)$ The test $d_{\lambda(\theta_1, \theta_0; \alpha)}$ given by

$$d_{\lambda(\theta_1, \theta_0; \alpha)}(y) = 0 \quad \text{iff} \quad y < \theta_0 + \Phi^{-1}(1 - \alpha)$$

can be implemented *without* explicit knowledge of θ_1 . By the discussion of Section 8.5, it is plain that a UMP test $d_{\text{UMP}}(\alpha)$ of size α exists with

$$C(d_{\text{UMP}}(\alpha)) = \{y \in \mathbb{R} : y < \theta_0 + \Phi^{-1}(1 - \alpha)\} = (-\infty, \theta_0 + \Phi^{-1}(1 - \alpha)).$$

Case II: $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = (-\infty, \theta_0)$ Here, the test $d_{\lambda(\theta_1, \theta_0; \alpha)}$ given by

$$d_{\lambda(\theta_1, \theta_0; \alpha)}(y) = 0 \quad \text{iff} \quad y > \theta_0 + \Phi^{-1}(\alpha).$$

can also be implemented without explicit knowledge of θ_1 . Again, the discussion of Section 8.5 shows that a UMP test $d_{\text{UMP}}(\alpha)$ of size α exists with

$$C(d_{\text{UMP}}(\alpha)) = \{y \in \mathbb{R} : y > \theta_0 + \Phi^{-1}(\alpha)\} = (\theta_0 + \Phi^{-1}(\alpha), \infty).$$

Case III: $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = \mathbb{R} - \{\theta_0\}$ It is plain from the discussion in Cases I and II that a UMP test $d_{\text{UMP}}(\alpha)$ does not exist.

Case IV: $\Theta_0 = (-\infty, \theta_0]$ and $\Theta_1 = (\theta_0, \infty)$ We begin by picking σ arbitrary in Θ_0 , so $\sigma \leq \theta_0$. Consider now the composite hypothesis testing problem

$$\begin{aligned} H_1 : Y &\sim N(\theta, 1), \quad \theta \in \Theta_1 \\ H_0 : Y &\sim N(\sigma, 1) \end{aligned} \quad (8.32)$$

The discussion in Case I also shows that a UMP test of size α exists; it was identified as the test $d^*(\alpha; \sigma) : \mathbb{R}^k \rightarrow \{0, 1\}$ given by

$$d^*(\alpha; \sigma)(y) = 0 \quad \text{iff} \quad y < \sigma + \Phi^{-1}(1 - \alpha).$$

This means that for every θ in Θ_1 , it holds that

$$\begin{aligned} \mathbb{P}_\theta [d^*(\alpha; \sigma)(Y) = 1] &= \mathbb{P}_\theta [Y \geq \sigma + \Phi^{-1}(1 - \alpha)] \\ &\geq \mathbb{P}_\theta [d(Y) = 1] \end{aligned} \quad (8.33)$$

for every test d in $\mathcal{D}_{\sigma, \alpha}$, i.e., for every test in \mathcal{D} such that $\mathbb{P}_\sigma [d(Y) = 1] \leq \alpha$.

Now recall that $\mathcal{D}_{\Theta_0, \alpha}$ is contained in $\mathcal{D}_{\sigma, \alpha}$ since

$$\mathcal{D}_{\Theta_0, \alpha} = \left\{ d \in \mathcal{D} : \sup_{\sigma' \leq \theta_0} \mathbb{P}_{\sigma'} [d(Y) = 1] \leq \alpha \right\}.$$

Therefore, for every θ in Θ_1 , we have

$$\mathbb{P}_\theta [d(Y) = 1] \leq \mathbb{P}_\theta [Y \geq \sigma + \Phi^{-1}(1 - \alpha)], \quad d \in \mathcal{D}_{\Theta_0, \alpha}$$

by virtue of (8.33) applied to the smaller class $\mathcal{D}_{\Theta_0, \alpha}$ of tests. Thus, with $\sigma \leq \theta_0$, the test $d^*(\alpha; \sigma)$ would be the UMP test $d_{\text{UMP}}(\alpha)$ we seek if only it belonged to $\mathcal{D}_{\Theta_0, \alpha}$.

Next we show that σ can be selected in Θ_0 such that the test $d^*(\alpha; \sigma)$ is indeed an element of $\mathcal{D}_{\Theta_0, \alpha}$, in which case it is the desired UMP test $d_{\text{UMP}}(\alpha)$. We shall prove that this happened if and only if $\sigma = \theta_0$: For arbitrary σ and σ' in \mathbb{R} we always have

$$\begin{aligned} F(\sigma, \sigma') &\equiv \mathbb{P}_{\sigma'} [Y \geq \sigma + \Phi^{-1}(1 - \alpha)] \\ &= \mathbb{P}_{\sigma'} [Y \geq \sigma' + (\sigma - \sigma') + \Phi^{-1}(1 - \alpha)] \\ &= \mathbb{P}_{\sigma'} [Y - \sigma' \geq (\sigma - \sigma') + \Phi^{-1}(1 - \alpha)] \\ &= 1 - \Phi((\sigma - \sigma') + \Phi^{-1}(1 - \alpha)). \end{aligned}$$

With σ in \mathbb{R} given, we see that the mapping $\mathbb{R} \rightarrow [0, 1] : \sigma' \rightarrow F(\sigma, \sigma')$ is *strictly increasing* with

$$F(\sigma, \sigma) = \alpha$$

since $\Phi(\Phi^{-1}(1 - \alpha)) = 1 - \alpha$, whence $F(\sigma, \sigma) < \alpha$ if $\sigma' < \sigma$ and $\alpha < F(\sigma, \sigma)$ if $\sigma < \sigma'$. It now follows that

$$\sup_{\sigma' \leq \theta_0} \mathbb{P}_{\sigma'} [d^*(\alpha; \sigma)(Y) = 1] > \alpha \quad \text{if } \sigma < \theta_0$$

while

$$\sup_{\sigma' \leq \theta_0} \mathbb{P}_{\sigma'} [d^*(\alpha; \theta_0)(Y) = 1] = \alpha.$$

This shows that the test $d^*(\alpha; \theta_0)$ is indeed an element of $\mathcal{D}_{\Theta_0, \alpha}$, and therefore can be used to implement the desired UMP test $d_{\text{UMP}}(\alpha)$. ■

8.9 Reformulating the Bayesian cost

Fix d in \mathcal{D} . By iterated conditioning we conclude that

$$\begin{aligned} J(d) &= \mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))] \\ &= \mathbb{E}[\mathbb{E}[\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))|\boldsymbol{\vartheta}] | H]] \end{aligned} \tag{8.34}$$

with

$$\begin{aligned} &\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))|\boldsymbol{\vartheta}] \\ &= C(\boldsymbol{\vartheta}, 0) \cdot \mathbb{P}[d(\mathbf{Y}) = 0|\boldsymbol{\vartheta}] + C(\boldsymbol{\vartheta}, 1) \cdot \mathbb{P}[d(\mathbf{Y}) = 1|\boldsymbol{\vartheta}] \\ &= C(\boldsymbol{\vartheta}, 0) \cdot \mathbb{P}[\mathbf{Y} \in C(d)|\boldsymbol{\vartheta}] + C(\boldsymbol{\vartheta}, 1) \cdot \mathbb{P}[\mathbf{Y} \notin C(d)|\boldsymbol{\vartheta}] \\ &= C(\boldsymbol{\vartheta}, 0) \int_{C(d)} f_{\boldsymbol{\vartheta}}(\mathbf{y}) dF(\mathbf{y}) + C(\boldsymbol{\vartheta}, 1) \int_{C(d)^c} f_{\boldsymbol{\vartheta}}(\mathbf{y}) dF(\mathbf{y}). \end{aligned} \tag{8.35}$$

On the other hand, for each $h = 0, 1$, we note as before that

$$\begin{aligned}
\mathbb{P}[\boldsymbol{\vartheta} \leq \mathbf{t} | H = h] &= \frac{\mathbb{P}[\boldsymbol{\vartheta} \leq \mathbf{t}, H = h]}{\mathbb{P}[H = h]} \\
&= \frac{\mathbb{P}[\boldsymbol{\vartheta} \leq \mathbf{t}, \boldsymbol{\vartheta} \in \Theta_h]}{\mathbb{P}[\boldsymbol{\vartheta} \in \Theta_h]} \\
&= \frac{\int_{(-\infty, \mathbf{t}] \cap \Theta_h} dK(\boldsymbol{\tau})}{\int_{\Theta_h} dK(\boldsymbol{\tau})} \\
&= \frac{\int_{-\infty}^{\mathbf{t}} \mathbf{1}[\boldsymbol{\tau} \in \Theta_h] dK(\boldsymbol{\tau})}{\int_{\Theta_h} dK(\boldsymbol{\tau})}, \quad \mathbf{t} \in \mathbb{R}^p. \quad (8.36)
\end{aligned}$$

Therefore, the conditional distribution of $\boldsymbol{\vartheta}$ given $H = h$ is absolutely continuous with respect to the probability distribution $K : \mathbb{R}^p \rightarrow [0, 1]$, the corresponding probability density function $k_h : \mathbb{R}^p \rightarrow \mathbb{R}_+$ being given by

$$k_h(\mathbf{t}) = \frac{1}{\int_{\Theta_h} dK(\boldsymbol{\tau})} \cdot \mathbf{1}[\mathbf{t} \in \Theta_h], \quad \mathbf{t} \in \mathbb{R}^p.$$

Next, we get

$$\begin{aligned}
&\mathbb{E}[\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y})) | \boldsymbol{\vartheta}] | H = h] \\
&= \mathbb{E}\left[C(\boldsymbol{\vartheta}, 0) \int_{C(d)} f_{\boldsymbol{\vartheta}}(\mathbf{y}) dF(\mathbf{y}) + C(\boldsymbol{\vartheta}, 1) \int_{C(d)^c} f_{\boldsymbol{\vartheta}}(\mathbf{y}) dF(\mathbf{y}) \Big| H = h\right] \\
&= \int_{\mathbb{R}^p} \left(C(\mathbf{t}, 0) \int_{C(d)} f_{\mathbf{t}}(\mathbf{y}) dF(\mathbf{y}) + C(\mathbf{t}, 1) \int_{C(d)^c} f_{\mathbf{t}}(\mathbf{y}) dF(\mathbf{y}) \right) k_h(\mathbf{t}) dK(\mathbf{t}) \\
&= \int_{C(d)} \left(\int_{\mathbb{R}^p} C(\mathbf{t}, 0) f_{\mathbf{t}}(\mathbf{y}) k_h(\mathbf{t}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \\
&\quad + \int_{C(d)^c} \left(\int_{\mathbb{R}^p} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) k_h(\mathbf{t}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \\
&= \int_{C(d)} \frac{\int_{\Theta_h} C(\mathbf{t}, 0) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_h} dK(\mathbf{t})} dF(\mathbf{y}) + \int_{C(d)^c} \frac{\int_{\Theta_h} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_h} dK(\mathbf{t})} dF(\mathbf{y}).
\end{aligned}$$

Next we observe that

$$\begin{aligned}
J(d) &= \mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))] \\
&= \mathbb{E}[\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y})) | \boldsymbol{\vartheta}] | H = 0] \mathbb{P}[H = 0] \\
&\quad + \mathbb{E}[\mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y})) | \boldsymbol{\vartheta}] | H = 1] \mathbb{P}[H = 1],
\end{aligned}$$

with elementary algebra showing that

$$\begin{aligned}
J(d) &= \mathbb{E}[C(\boldsymbol{\vartheta}, d(\mathbf{Y}))] \\
&= \left(\int_{C(d)} \frac{\int_{\Theta_0} C(\mathbf{t}, 0) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_0} dK(\mathbf{t})} dF(\mathbf{y}) \right) \cdot \mathbb{P}[H = 0] \\
&\quad + \left(\int_{C(d)^c} \frac{\int_{\Theta_0} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_0} dK(\mathbf{t})} dF(\mathbf{y}) \right) \cdot \mathbb{P}[H = 0] \\
&\quad + \left(\int_{C(d)} \frac{\int_{\Theta_1} C(\mathbf{t}, 0) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_1} dK(\mathbf{t})} dF(\mathbf{y}) \right) \cdot \mathbb{P}[H = 1] \\
&\quad + \left(\int_{C(d)^c} \frac{\int_{\Theta_1} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t})}{\int_{\Theta_1} dK(\mathbf{t})} dF(\mathbf{y}) \right) \cdot \mathbb{P}[H = 1] \\
&= \int_{C(d)} \left(\int_{\Theta} C(\mathbf{t}, 0) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \\
&\quad + \int_{C(d)^c} \left(\int_{\Theta} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \tag{8.37}
\end{aligned}$$

as we group like terms after noting that $\mathbb{P}[H = h] = \int_{\Theta_h} dK(\mathbf{t})$.

Finally we conclude that

$$\begin{aligned}
J(d) &= \int_{\mathbb{R}^k} \left(\int_{\Theta} C(\mathbf{t}, 1) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \\
&\quad + \int_{C(d)} \left(\int_{\Theta} (C(\mathbf{t}, 0) - C(\mathbf{t}, 1)) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) \right) dF(\mathbf{y}) \tag{8.38}
\end{aligned}$$

and the Bayesian test $d^* : \mathbb{R}^k \rightarrow \{0, 1\}$ is therefore given by

$$d^*(\mathbf{y}) = 0 \quad \text{iff} \quad \int_{\Theta} (C(\mathbf{t}, 0) - C(\mathbf{t}, 1)) f_{\mathbf{t}}(\mathbf{y}) dK(\mathbf{t}) < 0. \tag{8.39}$$

The Bayesian test d^* is *not* an element of the class of tests $\{d_{\eta}, \eta \geq 0\}$ introduced at (8.12). ■

8.10 Exercises

Part II

ESTIMATION THEORY

Part III
APENDICES

Chapter 9

Useful facts from Real Analysis

9.1 Limits in \mathbb{R}

We refer to any mapping $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ as a (\mathbb{R} -valued) *sequence*; sometimes we shall also use the notation $\{a_n, n = 1, 2, \dots\}$.

A sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ converges to a^* in \mathbb{R} if for every $\varepsilon > 0$, there exists an integer $n^*(\varepsilon)$ in \mathbb{N}_0 such that

$$|a_n - a^*| \leq \varepsilon, \quad n \geq n^*(\varepsilon). \quad (9.1)$$

We write $\lim_{n \rightarrow \infty} a_n = a^*$, and refer to the scalar a^* as the *limit* of the sequence.

Sometimes it is desirable to make sense of situations where the values of the sequence become either unboundedly large or unboundedly negative, in which case we shall write $\lim_{n \rightarrow \infty} a_n = \infty$ and $\lim_{n \rightarrow \infty} a_n = -\infty$, respectively. A precise definition of such occurrences is as follows: We write $\lim_{n \rightarrow \infty} a_n = \infty$ to signify that for every $M > 0$, there exists a integer $n^*(M)$ in \mathbb{N}_0 such that

$$a_n > M, \quad n \geq n^*(M). \quad (9.2)$$

It is now natural to define $\lim_{n \rightarrow \infty} a_n = -\infty$ whenever $\lim_{n \rightarrow \infty} (-a_n) = \infty$.

If there exists a^* in $\mathbb{R} \cup \{\pm\infty\}$ such that $\lim_{n \rightarrow \infty} a_n = a^*$, we shall simply say that the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ *converges* or *is convergent* (without any reference to its limit). Sometimes we shall also say that the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ converges (or is convergent) *in* \mathbb{R} to indicate that the limit a^* is an element of \mathbb{R} (thus finite).

Convergence is guaranteed under conditions of monotonicity: A sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ is *monotone non-decreasing* (resp. *non-increasing*) if $a_n \leq a_{n+1}$ (resp. $a_{n+1} \leq a_n$) for all $n = 1, 2, \dots$

Fact 9.1.1 *A monotone sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ is always convergent, although its limit may be $\pm\infty$.*

9.2 Accumulation points

Since not all sequences converge, it is important to understand how can non-convergence occur. To that end, consider a sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$. A *subsequence* of the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ is any sequence of the form $\mathbb{N}_0 \rightarrow \mathbb{R} : k \rightarrow a_{n_k}$ where

$$n_k < n_{k+1}, \quad k = 1, 2, \dots$$

This forces $\lim_{k \rightarrow \infty} n_k = \infty$.

An *accumulation point* for the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ is defined as any element a^* in $\mathbb{R} \cup \{\pm\infty\}$ such that

$$\lim_{k \rightarrow \infty} a_{n_k} = a^*$$

for *some* subsequence $\mathbb{N}_0 \rightarrow \mathbb{R} : k \rightarrow a_{n_k}$.

A convergent sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ has exactly *one* accumulation point, namely its limit. if the sequence *does not* converge, it must necessarily have distinct accumulation points, in which case there is a smallest and a largest accumulation point. The next definition formalizes this observation: Given a sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$, the quantities

$$\limsup_{n \rightarrow \infty} a_n = \inf_{n \geq 1} \left(\sup_{m \geq n} a_m \right)$$

and

$$\liminf_{n \rightarrow \infty} a_n = \sup_{n \geq 1} \left(\inf_{m \geq n} a_m \right)$$

are known as the *limsup* and *liminf* of the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$, respectively.

It is plain that

$$\inf_{m \geq n} a_m \leq \sup_{m \geq n} a_m, \quad n = 1, 2, \dots$$

and that the sequences $\{\inf_{m \geq n} a_m, n = 1, 2, \dots\}$ and $\{\sup_{m \geq n} a_m, n = 1, 2, \dots\}$ are non-decreasing and non-increasing, respectively, with

$$\lim_{n \rightarrow \infty} \left(\inf_{m \geq n} a_m \right) = \liminf_{n \rightarrow \infty} a_n$$

and

$$\lim_{n \rightarrow \infty} \left(\sup_{m \geq n} a_m \right) = \limsup_{n \rightarrow \infty} a_n,$$

hence the terminology. A useful characterization of convergence can now be provided in terms of the limsup and liminf.

Fact 9.2.1 Consider a sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$: If it converges to an element a^* (in $\mathbb{R} \cup \{\pm\infty\}$), then

$$\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n = a^*.$$

Conversely, if $\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n = a^*$ for some a^* in $\mathbb{R} \cup \{\pm\infty\}$, then the sequence $a : \mathbb{N}_0 \rightarrow \mathbb{R}$ converges to a^* .

9.3 Continuous functions

Let I denote a subset of \mathbb{R} . A function $g : I \rightarrow \mathbb{R}$ is said to be *left-continuous* at x in I if for any sequence $a : \mathbb{N}_0 \rightarrow I$ such that $\lim_{n \rightarrow \infty} a_n = x$ with $a_n < x$ for all $n = 1, 2, \dots$ we have

$$\lim_{n \rightarrow \infty} g(a_n) = g(x). \quad (9.3)$$

Similarly, a function $g : I \rightarrow \mathbb{R}$ is said to be *right-continuous* at x in I if for any sequence $a : \mathbb{N}_0 \rightarrow I$ such that $\lim_{n \rightarrow \infty} a_n = x$ with $x < a_n$ for all $n = 1, 2, \dots$ we have (9.3).

Finally, a function $g : I \rightarrow \mathbb{R}$ is said to be *continuous* at x in I if it is both left-continuous and right-continuous at x . This is equivalent to (9.3) holding for any sequence $a : \mathbb{N}_0 \rightarrow I$ such that $\lim_{n \rightarrow \infty} a_n = x$.

A function $g : I \rightarrow \mathbb{R}$ is left-continuous (resp. right-continuous, continuous) on I if the function $g : I \rightarrow \mathbb{R}$ is left-continuous (resp. right-continuous, continuous) at every point x in I .

9.4 Convex functions

Let I denote an interval in \mathbb{R} . A function $g : I \rightarrow \mathbb{R}$ is said to be *convex* if for every x_0 and x_1 in I , it holds that

$$g((1 - \lambda)x_0 + \lambda x_1) \leq (1 - \lambda)g(x_0) + \lambda g(x_1), \quad \lambda \in [0, 1]. \quad (9.4)$$

A function $g : I \rightarrow \mathbb{R}$ is said to be *concave* if the function $-g$ is convex. Here are some well-known facts concerning convex functions; the analog properties for concave functions are easily obtained *mutatis mutandis*.

Fact 9.4.1 *Let $g : I \rightarrow \mathbb{R}$ be a convex function. With $x < y < z$ in I , we have the basic inequalities*

$$\frac{g(y) - g(x)}{y - x} \leq \frac{g(z) - g(x)}{z - x} \quad (9.5)$$

and

$$\frac{g(z) - g(x)}{z - x} \leq \frac{g(z) - g(y)}{z - y}. \quad (9.6)$$

Proof. With $x < y < z$ in I , write

$$y = (1 - \lambda)x + \lambda z \quad \text{where} \quad \lambda = \frac{y - x}{z - x}.$$

It is plain that λ is an element of $(0, 1)$, and the convexity of g implies

$$\begin{aligned} g(y) &\leq (1 - \lambda)g(x) + \lambda g(z) \\ &= \frac{z - y}{z - x}g(x) + \frac{y - x}{z - x}g(z). \end{aligned} \quad (9.7)$$

Subtracting $g(x)$ from both sides of this inequality we get

$$\begin{aligned} g(y) - g(x) &\leq \frac{z - y}{z - x}g(x) + \frac{y - x}{z - x}g(z) - g(x) \\ &= \left(\frac{y - x}{z - x} \right) (g(z) - g(x)). \end{aligned}$$

and this establishes (9.5).

On the other hand, subtracting $g(z)$ from (9.7) we find that

$$\begin{aligned} g(y) - g(z) &\leq \frac{z-y}{z-x}g(x) + \frac{y-x}{z-x}g(z) - g(z) \\ &= \left(\frac{z-y}{z-x}\right)(g(x) - g(z)) \end{aligned}$$

and the inequality (9.5) readily follows. ■

Fact 9.4.2 *If the mapping $g : I \rightarrow \mathbb{R}$ is convex on an interval I , then it is also continuous on the interior of I .*

Continuity may fail at the boundary points as the next example shows: With $I = [0, 1]$, take the mapping $g : I \rightarrow \mathbb{R}$ given by

$$g(x) = \begin{cases} 0 & \text{if } x = 0 \\ x & \text{if } 0 < x \leq 1. \end{cases}$$

This function is clearly convex on $[0, 1]$ but it fails to be continuous at $x = 0$.

Proof. Pick x in the interior of I so that $[x - \tau, x + \tau] \subseteq I$ for some $\tau > 0$. With $0 < t < 1$ we note that $x + t\tau = t(x + \tau) + (1 - t)x$, whence $g(x + t\tau) \leq tg(x + \tau) + (1 - t)g(x)$ by the convexity of g and we conclude

$$\frac{g(x + t\tau) - g(x)}{t} \leq g(x + \tau) - g(x).$$

Similarly, since $x - t\tau = t(x - \tau) + (1 - t)x$, we have $g(x - t\tau) \leq tg(x - \tau) + (1 - t)g(x)$ by the convexity of g , so that

$$g(x) - g(x - \tau) \leq \frac{g(x) - g(x - t\tau)}{t}.$$

But $\frac{1}{2}(x - t\tau) + \frac{1}{2}(x + t\tau) = x$, and using convexity again leads to

$$g(x) - g(x - t\tau) \leq g(x + t\tau) - g(x),$$

whence

$$g(x) - g(x - \tau) \leq \frac{g(x - t\tau) - g(x)}{-t} \leq \frac{g(x + t\tau) - g(x)}{t} \leq g(x + \tau) - g(x).$$

Therefore, with

$$M \equiv \max(|g(x) - g(x - \tau)|, |g(x + \tau) - g(x)|),$$

we get

$$-M \leq \frac{g(x) - g(x - t\tau)}{-t} \leq \frac{g(x + t\tau) - g(x)}{t} \leq M, \quad 0 < t < 1.$$

Changing notation we see that

$$\left| \frac{g(x + h) - g(x)}{h} \right| \leq M, \quad |h| \leq \tau \quad (9.8)$$

and the continuity of g at x is now immediate since $\lim_{n \rightarrow \infty} g(a_n) = g(x)$ for any sequence I -valued sequence $a : \mathbb{N}_0 \rightarrow I$ such that $\lim_{n \rightarrow \infty} a_n = x$. ■

In the course of the proof of Fact 9.4.2 we saw that the property (9.8), known as local Lipschitz continuity, holds.

Fact 9.4.3 *If the mapping $g : I \rightarrow \mathbb{R}$ is convex on an interval I , then it is also locally Lipschitz at every point in the interior of I .*

This paves the way for the following differentiability result.

Fact 9.4.4 *If the mapping $g : I \rightarrow \mathbb{R}$ is convex on some interval I , then its left and right-derivatives always exist at every point of continuity in I (and whenever appropriate at the boundary points of I).*

Proof. Pick ξ in the interior of I so that $(\xi - \tau, \xi + \tau) \subseteq I$ for some $\tau > 0$. By virtue of (9.5) we see that the mapping $t \rightarrow \frac{g(\xi+t) - g(\xi)}{t}$ is monotone increasing on

the interval $(0, \tau)$ – Just use $x = \xi$, $y = \xi + t_1$ and $z = \xi + t_2$ with $0 < t_1 < t_2 < \tau$. It follows that the limit defining the right-derivative at ξ , namely

$$\frac{d^+}{dx}g(\xi) = \lim_{t \downarrow 0} \frac{g(\xi + t) - g(\xi)}{t},$$

always exists.

Similarly, by virtue of (9.6) we see that the mapping $t \rightarrow \frac{g(\xi - t) - g(\xi)}{-t}$ is monotone increasing on the interval $(0, \tau)$ – Just use $z = \xi$, $y = \xi - t_1$ and $x = \xi - t_2$ with $0 < t_1 < t_2 < \tau$. It follows that the limit defining the left-derivative at ξ , namely

$$\frac{d^-}{dx}g(\xi) = \lim_{t \downarrow 0} \frac{g(\xi - t) - g(\xi)}{-t},$$

always exists. ■

9.5 Measurable spaces

Let S denote an arbitrary non-empty set. A non-empty collection \mathcal{S} of subsets of S is a σ -field (also known as a σ -algebra) on S if

- (i) \mathcal{S} contains the empty set \emptyset ;
- (ii) \mathcal{S} is closed under complementarity: If $E \in \mathcal{S}$, then $E^c \in \mathcal{S}$ (where E^c is the complement of E in S); and
- (iii) \mathcal{S} is closed under countable union: With I a countable index set, if $E_i \in \mathcal{S}$ for each $i \in I$, then $\cup_{i \in I} E_i \in \mathcal{S}$.

The pair (S, \mathcal{S}) is sometimes referred to as a *measurable space*. For every non-empty set S , there are at least two distinct σ -fields on S , namely the trivial σ -field $\mathcal{S}_{\text{Triv}} = \{\emptyset, S\}$ and the complete σ -field $\mathcal{P}(S)$ (where $\mathcal{P}(S)$ denotes the power set of S).

If \mathcal{S}_1 and \mathcal{S}_2 are two σ -fields on S , we say that \mathcal{S}_1 contains \mathcal{S}_2 , written $\mathcal{S}_1 \subseteq \mathcal{S}_2$, if any element of \mathcal{S}_1 is an element of \mathcal{S}_2 . Thus, for any σ -field \mathcal{S} , we have $\mathcal{S}_{\text{Triv}} \subseteq \mathcal{S} \subseteq \mathcal{P}(S)$.

If \mathcal{G} is a collection of subsets of S , then $\sigma(\mathcal{G})$ is defined as the *smallest* σ -field on S which contains \mathcal{G} , i.e., every element of \mathcal{G} is also an element of $\sigma(\mathcal{G})$. We shall refer to $\sigma(\mathcal{G})$ as the σ -field on S generated by \mathcal{G} .

9.6 Borel measurability

With A denoting a subset of \mathbb{R}^p for some positive integer p , we define the σ -field $\mathcal{B}(A)$ to be

$$\mathcal{B}(A) \equiv \sigma(\mathcal{O}(A))$$

where $\mathcal{O}(A)$ denotes the collection of all *open sets* contained in A . In particular,

$$\mathcal{B}(\mathbb{R}^p) \equiv \sigma(\mathcal{O}(\mathbb{R}^p))$$

where $\mathcal{O}(\mathbb{R}^p)$ denotes the collection of all *open sets* contained in \mathbb{R}^p .

Consider an arbitrary set S equipped with a σ -field \mathcal{S} . A mapping $g : S \rightarrow \mathbb{R}^p$ is said to be a *Borel mapping* if the conditions

$$g^{-1}(B) \in \mathcal{S}, \quad B \in \mathcal{B}(\mathbb{R}^p) \tag{9.9}$$

are all satisfied where

$$g^{-1}(B) \equiv \{s \in S : g(s) \in B\}.$$

Fact 9.6.1 Let \mathcal{G} denote a collection of subsets of \mathbb{R}^p which generates the Borel σ -field $\mathcal{B}(\mathbb{R}^p)$, i.e.,

$$\mathcal{B}(\mathbb{R}^p) = \sigma(\mathcal{G}). \tag{9.10}$$

It holds that the mapping $g : S \rightarrow \mathbb{R}^p$ is a Borel mapping if and only if the weaker set of conditions

$$g^{-1}(E) \in \mathcal{S}, \quad E \in \mathcal{G} \tag{9.11}$$

holds.

There are many generators known for the Borel σ -field $\mathcal{B}(\mathbb{R}^p)$. For instance, we have (9.10) with

- $\mathcal{G} = \mathcal{R}_{\text{open}}(\mathbb{R}^p)$ where $\mathcal{R}_{\text{open}}(\mathbb{R}^p)$ is the collection of all *finite open rectangles*, i.e.,

$$\mathcal{R}_{\text{open}}(\mathbb{R}^p) \equiv \left\{ I_1 \times \dots \times I_p, \quad \begin{array}{l} I_k \in \mathcal{I}(\mathbb{R}) \\ k = 1, \dots, p \end{array} \right\}$$

where

$$\mathcal{I}(\mathbb{R}) = \{(a, b) : a, b \in \mathbb{R}\}$$

- $\mathcal{G} = \mathcal{R}_{\text{SW}}(\mathbb{R}^p)$ where $\mathcal{R}_{\text{SW}}(\mathbb{R}^p)$ is the collection of all *closed Southwest rectangles*, i.e.,

$$\mathcal{R}_{\text{SW}}(\mathbb{R}^p) \equiv \left\{ I_1 \times \dots \times I_p, \begin{array}{l} I_k = (-\infty, a_k] \\ a_k \in \mathbb{R} \\ k = 1, \dots, p \end{array} \right\}.$$

It follows from the discussion above that a mapping $g : S \rightarrow \mathbb{R}^p$ is a Borel mapping if the seemingly weaker conditions

$$\left\{ s \in S : g(s) \in \prod_{i=1}^p (-\infty, a_i] \right\} \in \mathcal{S}, \quad (a_1, \dots, a_p) \in \mathbb{R}^p$$

all hold. Equivalently, a mapping $g : S \rightarrow \mathbb{R}^p$ is a Borel mapping if

$$\{s \in S : g_k(s) \leq a_k, k = 1, \dots, p\} \in \mathcal{S}, \quad (a_1, \dots, a_p) \in \mathbb{R}^p$$

where it is understood that

$$g(s) = (g_1(s), \dots, g_p(s)), \quad s \in S.$$

It is now plain that for each $k = 1, \dots, p$, the component mapping $g_k : S \rightarrow \mathbb{R}$ is also a Borel mapping – Just take $a_\ell = \infty$ for all $\ell = 1, \dots, k$ different from k . Conversely, since

$$\{s \in S : g_k(s) \leq a_k, k = 1, \dots, p\} = \bigcap_{k=1}^p \{s \in S : g_k(s) \leq a_k\}$$

for arbitrary (a_1, \dots, a_p) in \mathbb{R}^p , we see that the mapping $g : S \rightarrow \mathbb{R}^p$ is a Borel mapping if and only if each of the component mappings $g_1 : S \rightarrow \mathbb{R}, \dots, g_p : S \rightarrow \mathbb{R}$ is a Borel mapping.

Most (if not all) mappings $\mathbb{R}^p \rightarrow \mathbb{R}^q$ encountered in applications are Borel mappings. Furthermore, any *continuous* mapping $\mathbb{R}^p \rightarrow \mathbb{R}^q$ can be shown to be a Borel mapping!

Chapter 10

Useful facts from Probability Theory

10.1 Probability models

Probabilistic reasoning assumes the availability of a probability triple $(\Omega, \mathcal{F}, \mathbb{P})$ where: (i) The *sample space* Ω is the collection of all outcomes (samples) generated by the random experiment \mathcal{E} of interest; (ii) Events are collections of outcomes, and the collection of events whose likelihood of occurrence can be defined is a σ -field \mathcal{F} on Ω ; and (iii) The “likelihood” of occurrence to events in \mathcal{F} is assigned through a *probability measure* \mathbb{P} defined on \mathcal{F} .

With Ω an arbitrary set, a non-empty collection \mathcal{F} of subsets of Ω is a σ -field (also known as a σ -algebra) on Ω if \mathcal{F} (i) contains the empty set \emptyset ; (ii) is closed under complementarity: If $E \in \mathcal{F}$, then $E^c \in \mathcal{F}$ (iii) is closed under countable union: With I a *countable* index set, if $E_i \in \mathcal{F}$ for each $i \in I$, then $\cup_{i \in I} E_i \in \mathcal{F}$.

The “likelihood” of occurrence to events in \mathcal{F} is assigned through a *probability measure* \mathbb{P} defined on \mathcal{F} .

A *probability (measure)* \mathbb{P} on the σ -field \mathcal{F} (or on (Ω, \mathcal{F})) is a mapping $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ such that (i) $\mathbb{P}[\emptyset] = 0$ and $\mathbb{P}[\Omega] = 1$; (ii) σ -additivity: With I a countable index set, if $E_i \in \mathcal{F}$ for each $i \in I$, then

$$\mathbb{P}[\cup_{i \in I} E_i] = \sum_{i \in I} \mathbb{P}[E_i]$$

whenever the subsets $\{E_i, i \in I\}$ are *pairwise disjoint*, namely

$$E_i \cap E_j = \emptyset, \quad \begin{array}{l} i \neq j \\ i, j \in I \end{array}$$

10.2 Random variables

All random variables (rvs) can always be thought as being defined on some given probability triple $(\Omega, \mathcal{F}, \mathbb{P})$ where Ω is the sample space, \mathcal{F} is a σ -field of events on Ω and \mathbb{P} is a probability measure on \mathcal{F} .

Given a probability triple $(\Omega, \mathcal{F}, \mathbb{P})$, a mapping $X : \Omega \rightarrow \mathbb{R}^p$ is an (\mathbb{R}^p -valued) *random variable* (rv) if

$$X^{-1}(B) = \{\omega \in \Omega : X(\omega) \in B\} \in \mathcal{F}, \quad B \in \mathcal{B}(\mathbb{R}^p).$$

In other words, the mapping $X : \Omega \rightarrow \mathbb{R}^p$ is a rv if it is a Borel mapping $X : \Omega \rightarrow \mathbb{R}^p$ – Here $S = \Omega$ and $\mathcal{S} = \mathcal{F}$. We shall often write $[X \in B]$ in lieu of $X^{-1}(B)$ and $\mathbb{P}[X \in B]$ for $\mathbb{P}[[X \in B]]$.

The *probability distribution (function)* $F_X : \mathbb{R}^p \rightarrow [0, 1]$ of the rv X is defined by

$$\begin{aligned} F_X(x) &\equiv \mathbb{P}[X \in (-\infty, x_1] \times \dots \times (-\infty, x_p)] \\ &= \mathbb{P}[X_1 \leq x_1, \dots, X_p \leq x_p], \quad x = (x_1, \dots, x_p) \in \mathbb{R}^p. \end{aligned} \quad (10.1)$$

with the notation $X = (X_1, \dots, X_p)$.

It turns out that there is as much probabilistic information in the probability distribution $F_X : \mathbb{R}^p \rightarrow [0, 1]$ of the rv X as in

$$\{\mathbb{P}[X \in B], B \in \mathcal{B}(\mathbb{R}^p)\}$$

In fact, knowledge of $F_X : \mathbb{R}^p \rightarrow \mathbb{R}$ allows a *unique* reconstruction of

$$\{\mathbb{P}[X \in B], B \in \mathcal{B}(\mathbb{R}^p)\}.$$

10.3 Probability distributions

Properties of F_X (Case $p = 1$): It is easy to see that the following properties hold:

- Monotonicity:

$$F_X(x) \leq F_X(y), \quad x, y \in \mathbb{R}$$

- Right-continuous:

$$\lim_{y \downarrow x} F_X(y) = F_X(x), \quad x \in \mathbb{R}$$

- Left limit exists:

$$\lim_{y \uparrow x} F_X(y) = F_X(x-) \quad \text{with} \quad \mathbb{P}[X = x] = F_X(x) - F_X(x-), \quad x \in \mathbb{R}$$

- Behavior at infinity: Monotonically

$$\lim_{x \rightarrow -\infty} F_X(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F_X(x) = 1$$

A *probability distribution (function)* on \mathbb{R} is any mapping $F : \mathbb{R} \rightarrow [0, 1]$ such that

- Monotonicity:

$$F(x) \leq F(y), \quad x, y \in \mathbb{R}$$

- Right-continuous:

$$\lim_{y \downarrow x} F(y) = F(x), \quad x \in \mathbb{R}$$

- Left limit exists:

$$\lim_{y \uparrow x} F(y) = F(x-) \quad x \in \mathbb{R}$$

- Behavior at infinity: Monotonically

$$\lim_{x \rightarrow -\infty} F(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F(x) = 1$$

Any rv $X : \Omega \rightarrow \mathbb{R}$ generates a probability distribution function $F_X : \mathbb{R} \rightarrow [0, 1]$. Conversely, for any probability distribution function $F : \mathbb{R} \rightarrow [0, 1]$, there exists a probability triple $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$ and a rv $\tilde{X} : \tilde{\Omega} \rightarrow \mathbb{R}$ defined on it such that

$$\tilde{\mathbb{P}}[\tilde{X} \leq x] = F(x), \quad x \in \mathbb{R}$$

This is the basis of Monte-Carlo simulation. There exists a multi-dimensional analog to this fact.

Proof. Take $\tilde{\Omega} = [0, 1]$, $\tilde{\mathcal{F}} = \mathcal{B}([0, 1])$ and $\tilde{\mathbb{P}} = \lambda$. Define the rv $\tilde{X} : \tilde{\Omega} \rightarrow \mathbb{R}$ by setting

$$\tilde{X}(\tilde{\omega}) = F^{-}(\tilde{\omega}), \quad \tilde{\omega} \in [0, 1]$$

where $F^{-} : [0, 1] \rightarrow [-\infty, \infty]$ is the *generalized inverse* of F given by

$$F^{-}(u) = \inf\{x \in \mathbb{R} : u \leq F(x)\}, \quad 0 \leq u \leq 1.$$

with the understanding that $F^{-}(u) = \infty$ if the defining set is empty, i.e., $F(x) < u$ for all x in \mathbb{R} .

Discrete distributions

A rv $X : \Omega \rightarrow \mathbb{R}^p$ is a *discrete* rv if there exists a *countable* subset $S \subseteq \mathbb{R}^p$ such that

$$\mathbb{P}[X \in S] = 1.$$

Note that

$$\mathbb{P}[X \in B] = \sum_{x \in S \cap B} \mathbb{P}[X = x], \quad B \in \mathcal{B}(\mathbb{R}^p).$$

It is often more convenient to characterize the distributional properties of the rv X through its probability mass function (pmf) of the rv X given by

$$\mathbf{p}_X \equiv (p_X(x), x \in S)$$

with

$$p_X(x) = \mathbb{P}[X = x], \quad x \in S.$$

Absolutely continuous distributions

A rv $X : \Omega \rightarrow \mathbb{R}^p$ is an (absolutely) continuous rv if there exists a Borel mapping $f_X : \mathbb{R}^p \rightarrow \mathbb{R}_+$ such that

$$\mathbb{P}[X_i \leq x_i, i = 1, \dots, p] = \int_{-\infty}^x f_X(\xi) d\xi, \quad x = (x_1, \dots, x_p) \in \mathbb{R}^p.$$

Properties of F_X when $p \geq 1$

- Monotonicity needs to be modified and now reads

$$\mathbb{P}[x_k < X_k \leq y_k] \geq 0, \quad \begin{array}{l} x_k < y_k \\ x_k, y_k \in \mathbb{R} \\ k = 1, \dots, p \end{array}$$

with the understanding that the quantity $\mathbb{P}[x_k < X_k \leq y_k]$ is expressed solely in terms of $F_X : \mathbb{R}^p \rightarrow [0, 1]$.

- Right-continuous:

$$\lim_{y \downarrow x} F_X(y) = F_X(x), \quad x \in \mathbb{R}^p$$

with the understanding that $y_k \downarrow x_k$ for each $k = 1, \dots, p$.

- Left limit exists:

$$\lim_{y \uparrow x} F_X(y) = F_X(x-) \quad \text{with} \quad \mathbb{P}[X = x] = F_X(y) - F_X(x-), \quad x \in \mathbb{R}^p$$

with the understanding that $y_k \uparrow x_k$ for each $k = 1, \dots, p$.

- Behavior at infinity:

$$\lim_{\min(x_k, k=1, \dots, p) \rightarrow -\infty} F_X(x) = 0$$

and

$$\lim_{\min(x_k, k=1, \dots, p) \rightarrow \infty} F_X(x) = 1$$

Independence of rvs

Consider a collection of rvs $\{X_i, i \in I\}$ which are all defined on some probability triple $(\Omega, \mathcal{F}, \mathbb{P})$. Assume that for each i in I , the rv X_i is a \mathbb{R}^{p_i} -valued rv for some positive integer p_i .

With I finite, we shall say that the rvs $\{X_i, i \in I\}$ are *mutually independent* if for each selection of B_i in $\mathcal{B}(\mathbb{R}^{p_i})$ for each i in I , the events

$$\{[X_i \in B_i], i \in I\}$$

are mutually independent. It is easy to see that this is equivalent to requiring

$$\mathbb{P}[\cap_{i \in I} [X_i \in B_i]] = \prod_{i \in I} \mathbb{P}[X_i \in B_i], \quad \begin{array}{l} B_i \in \mathcal{B}(R^{p_i}) \\ i \in I. \end{array}$$

More generally, with I arbitrary (and possibly uncountable), the rvs $\{X_i, i \in I\}$ are said to be mutually independent if for every finite subset $J \subseteq I$, the rvs $\{X_j, j \in J\}$ are mutually independent!

Product spaces

Some facts: Consider two arbitrary sets Ω_a and Ω_b (possibly identical). Let \mathcal{A} and \mathcal{B} denote non-empty collections of subsets of Ω_a and Ω_b , respectively. While the collection $\mathcal{A} \times \mathcal{B}$ is usually *not* a σ -field on $\Omega_a \times \Omega_b$, even when \mathcal{A} and \mathcal{B} are themselves σ -fields, it can be shown that

$$\sigma(\mathcal{A} \times \mathcal{B}) = \sigma(\sigma(\mathcal{A}) \times \sigma(\mathcal{B})).$$

10.4 Gaussian rvs

With scalar m and $\sigma^2 > 0$, the rv X is said to be a Gaussian rv with mean m and variance σ^2 , written $X \sim N(m, \sigma^2)$, if its cumulative probability distribution function is given by

$$\mathbb{P}[X \leq x] = \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{\xi-m}{\sigma}\right)^2} d\xi, \quad x \in \mathbb{R}.$$

The zero mean unit variance Gaussian rv is often referred to as a standard Gaussian rv; its probability density function $\varphi : \mathbb{R} \rightarrow \mathbb{R}_+$ is given by

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad x \in \mathbb{R}, \quad (10.2)$$

and its cumulative probability distribution function is then

$$\Phi(x) = \int_{-\infty}^x \varphi(t) dt, \quad x \in \mathbb{R}. \quad (10.3)$$

Obviously, if $X \sim N(m, \sigma^2)$ and $Z \sim N(0, 1)$, then X and $m + \sigma Z$ have the same distribution.

In the context of digital communications, it is customary to use the Q -function $Q : \mathbb{R}_+ \rightarrow [0, 1]$ given by

$$Q(x) = 1 - \Phi(x) = \int_x^{\infty} \varphi(t) dt, \quad x \geq 0.$$

Chapter 11

The classical limit theorems

The setting of the next four sections is as follows: The rvs $\{X_n, n = 1, 2, \dots\}$ are rvs defined on some probability triple $(\Omega, \mathcal{F}, \mathbb{P})$. With this sequence we associate the sums

$$S_n = \sum_{k=1}^n X_k, \quad n = 1, 2, \dots$$

Two types of results will be discussed: The first class of results are known as Laws of Large Numbers; they deal with the convergence of the sample averages

$$\bar{S}_n = \frac{1}{n} \sum_{k=1}^n X_k, \quad n = 1, 2, \dots$$

The second class of results are called Central Limit Theorems and provide a rate of convergence in the Laws Large Numbers.

11.1 Weak Laws of Large Numbers (I)

Laws of Large Numbers come in two types which are distinguished by the mode of convergence used. When convergence in probability is used, we refer to such results as weak Laws of Large Numbers. The most basic such results is given first.

Theorem 11.1.1 *Assume the rvs $\{X, X_n, n = 1, 2, \dots\}$ to be i.i.d. rvs with $\mathbb{E}[|X|^2] < \infty$. Then,*

$$\frac{S_n}{n} \xrightarrow{P} \mathbb{E}[X]. \quad (11.1)$$

As we now show, the finiteness of the second moment of X can be dropped.

Theorem 11.1.2 *Assume the rvs $\{X, X_n, n = 1, 2, \dots\}$ to be i.i.d. rvs with $\mathbb{E}[|X|] < \infty$. Then, we have*

$$\frac{S_n}{n} \xrightarrow{P} \mathbb{E}[X]. \quad (11.2)$$

11.2 The Strong Law of Large Numbers

Strong Laws of Large Numbers are given as convergence statements in the a.s. sense. The classical Strong Law of Large Numbers is given next.

Theorem 11.2.1 *Assume the rvs $\{X, X_n, n = 1, 2, \dots\}$ to be i.i.d. rvs with $\mathbb{E}[|X|] < \infty$. Then,*

$$\lim_{n \rightarrow \infty} \frac{S_n}{n} = \mathbb{E}[X] \quad a.s. \quad (11.3)$$

11.3 The Central Limit Theorem

The Central Limit Theorem completes the Law of Large Numbers, in that it provides some indication as to the rate at which convergence takes place.

Theorem 11.3.1 *Assume the rvs $\{X, X_n, n = 1, 2, \dots\}$ to be i.i.d. rvs with $\mathbb{E}[|X|^2] < \infty$. Then, we have*

$$\sqrt{n} \left(\frac{S_n}{n} - \mathbb{E}[X] \right) \Rightarrow_n \sqrt{\text{Var}[X]} \cdot U \quad (11.4)$$

where U is standard zero-mean unit-variance Gaussian rv.