

Mathematical Statistics

X_1, \dots, X_n iid $UV(\mu, \sigma^2)$. $\mu = \mu_0$ known.

$p_\sigma(x_1, \dots, x_n) = \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2 \right] \cdot \frac{1}{\sigma^n}$

$c(x) = \sum_{i=1}^n (x_i - \mu_0)^2 \quad c(\sigma^2) = -\frac{1}{2\sigma^2} \uparrow \text{ in } \sigma^2$

X_1, \dots, X_n iid Bernoulli(θ).

$p_\theta(x) = \theta^x (1-\theta)^{1-x} = \left(\frac{\theta}{1-\theta} \right)^x (1-\theta)$

$c(\theta) = \log \left[\frac{\theta}{1-\theta} \right] \uparrow \text{ in } \theta.$

The

Sara van de Geer

September 2010

Contents

1	Introduction	7
1.1	Some notation and model assumptions	7
1.2	Estimation	10
1.3	Comparison of estimators: risk functions	12
1.4	Comparison of estimators: sensitivity	12
1.5	Confidence intervals	13
1.5.1	Equivalence confidence sets and tests	13
1.6	Intermezzo: quantile functions	14
1.7	How to construct tests and confidence sets	14
1.8	An illustration: the two-sample problem	16
1.8.1	Assuming normality	17
1.8.2	A nonparametric test	18
1.8.3	Comparison of Student's test and Wilcoxon's test	20
1.9	How to construct estimators	21
1.9.1	Plug-in estimators	21
1.9.2	The method of moments	22
1.9.3	Likelihood methods	23
2	Decision theory	29
2.1	Decisions and their risk	29
2.2	Admissibility	31
2.3	Minimaxity	33
2.4	Bayes decisions	34
2.5	Intermezzo: conditional distributions	35
2.6	Bayes methods	36
2.7	Discussion of Bayesian approach (<i>to be written</i>)	39
2.8	Integrating parameters out (<i>to be written</i>)	39
2.9	Intermezzo: some distribution theory	39
2.9.1	The multinomial distribution	39
2.9.2	The Poisson distribution	41
2.9.3	The distribution of the maximum of two random variables	42
2.10	Sufficiency	42
2.10.1	Rao-Blackwell	44
2.10.2	Factorization Theorem of Neyman	45
2.10.3	Exponential families	47
2.10.4	Canonical form of an exponential family	48

2.10.5	Minimal sufficiency	53
3	Unbiased estimators	55
3.1	What is an unbiased estimator?	55
3.2	UMVU estimators	56
3.2.1	Complete statistics	59
3.3	The Cramer-Rao lower bound	62
3.4	Higher-dimensional extensions	66
3.5	Uniformly most powerful tests	68
3.5.1	An example	68
3.5.2	UMP tests and exponential families	71
3.5.3	Unbiased tests	74
3.5.4	Conditional tests	77
4	Equivariant statistics	81
4.1	Equivariance in the location model	81
4.2	Equivariance in the location-scale model (<i>to be written</i>)	86
5	Proving admissibility and minimaxity	87
5.1	Minimaxity	88
5.2	Admissibility	89
5.3	Inadmissibility in higher-dimensional settings (<i>to be written</i>)	95
6	Asymptotic theory	97
6.1	Types of convergence	97
6.1.1	Stochastic order symbols	99
6.1.2	Some implications of convergence	99
6.2	Consistency and asymptotic normality	101
6.2.1	Asymptotic linearity	101
6.2.2	The δ -technique	102
6.3	M-estimators	104
6.3.1	Consistency of M-estimators	106
6.3.2	Asymptotic normality of M-estimators	109
6.4	Plug-in estimators	114
6.4.1	Consistency of plug-in estimators	117
6.4.2	Asymptotic normality of plug-in estimators	118
6.5	Asymptotic relative efficiency	121
6.6	Asymptotic Cramer Rao lower bound	123
6.6.1	Le Cam's 3 rd Lemma	126
6.7	Asymptotic confidence intervals and tests	129
6.7.1	Maximum likelihood	131
6.7.2	Likelihood ratio tests	135
6.8	Complexity regularization (<i>to be written</i>)	139
7	Literature	141

These notes in English will closely follow *Mathematische Statistik*, by H.R. Künsch (2005), but are as yet incomplete. *Mathematische Statistik* can be used as supplementary reading material in German.

Mathematical rigor and clarity often bite each other. At some places, not all subtleties are fully presented. A snake will indicate this.



Chapter 1

Introduction

Statistics is about the mathematical modeling of observable phenomena, using stochastic models, and about analyzing data: estimating parameters of the model and testing hypotheses. In these notes, we study various estimation and testing procedures. We consider their theoretical properties and we investigate various notions of optimality.

1.1 Some notation and model assumptions

The data consist of measurements (observations) x_1, \dots, x_n , which are regarded as realizations of random variables X_1, \dots, X_n . In most of the notes, the X_i are real-valued: $X_i \in \mathbb{R}$ (for $i = 1, \dots, n$), although we will also consider some extensions to vector-valued observations.

Example 1.1.1 Fizeau and Foucault developed methods for estimating the speed of light (1849, 1850), which were later improved by Newcomb and Michelson. The main idea is to pass light from a rapidly rotating mirror to a fixed mirror and back to the rotating mirror. An estimate of the velocity of light is obtained, taking into account the speed of the rotating mirror, the distance travelled, and the displacement of the light as it returns to the rotating mirror.



Fig. 1

The data are Newcomb's measurements of the passage time it took light to travel from his lab, to a mirror on the Washington Monument, and back to his lab.

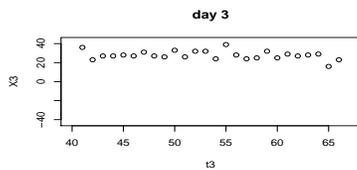
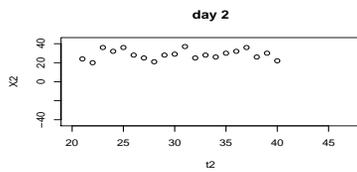
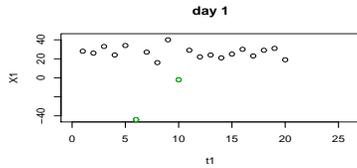
distance: 7.44373 km.

66 measurements on 3 consecutive days

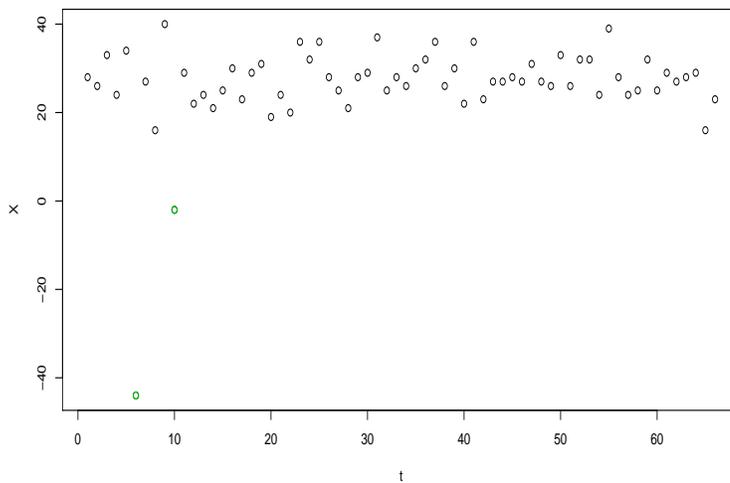
first measurement: 0.000024828 seconds= 24828 nanoseconds

The dataset has the deviations from 24800 nanoseconds.

The measurements on 3 different days:



All measurements in one plot:



One may estimate the speed of light using e.g. the mean, or the median, or Huber's estimate (see below). This gives the following results (for the 3 days separately, and for the three days combined):

	Day 1	Day 2	Day 3	All
Mean	21.75	28.55	27.85	26.21
Median	25.5	28	27	27
Huber	25.65	28.40	27.71	27.28

Table 1

The question which estimate is “the best one” is one of the topics of these notes.

Notation

The collection of observations will be denoted by $\mathbf{X} = \{X_1, \dots, X_n\}$. The distribution of \mathbf{X} , denoted by \mathbb{P} , is generally unknown. A statistical model is a collection of assumptions about this unknown distribution.

We will usually assume that the observations X_1, \dots, X_n are independent and identically distributed (i.i.d.). Or, to formulate it differently, X_1, \dots, X_n are i.i.d. copies from some population random variable, which we denote by X . The common distribution, that is: the distribution of X , is denoted by P . For $X \in \mathbb{R}$, the distribution function of X is written as

$$F(\cdot) = P(X \leq \cdot).$$

Recall that the distribution function F determines the distribution P (and vice versa).

Further model assumptions then concern the modeling of P . We write such a model as $P \in \mathcal{P}$, where \mathcal{P} is a given collection of probability measures, the so-called model class.

The following example will serve to illustrate the concepts that are to follow.

Example 1.1.2 Let X be real-valued. The location model is

$$\mathcal{P} := \{P_{\mu, F_0}(X \leq \cdot) := F_0(\cdot - \mu), \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0\}, \quad (1.1)$$

where \mathcal{F}_0 is a given collection of distribution functions. Assuming the expectation exist, we center the distributions in \mathcal{F}_0 to have mean zero. Then P_{μ, F_0} has mean μ . We call μ a location parameter. Often, only μ is the parameter of interest, and F_0 is a so-called nuisance parameter.

The class \mathcal{F}_0 is for example modeled as the class of all symmetric distributions, that is

$$\mathcal{F}_0 := \{F_0(x) = 1 - F_0(-x), \forall x\}. \quad (1.2)$$

This is an infinite-dimensional collection: it is not parametrized by a finite dimensional parameter. We then call F_0 an infinite-dimensional parameter.

A finite-dimensional model is for example

$$\mathcal{F}_0 := \{\Phi(\cdot/\sigma) : \sigma > 0\}, \quad (1.3)$$

where Φ is the standard normal distribution function.

Thus, the location model is

$$X_i = \mu + \epsilon_i, \quad i = 1, \dots, n,$$

with $\epsilon_1, \dots, \epsilon_n$ i.i.d. and, under model (1.2), symmetrically but otherwise unknown distributed and, under model (1.3), $\mathcal{N}(0, \sigma^2)$ -distributed with unknown variance σ^2 .

1.2 Estimation

A parameter is an aspect of the unknown distribution. An estimator T is some given function $T(\mathbf{X})$ of the observations \mathbf{X} . The estimator is constructed to estimate some unknown parameter, γ say.

In Example 1.1.2, one may consider the following estimators $\hat{\mu}$ of μ :

- The average

$$\hat{\mu}_1 := \frac{1}{n} \sum_{i=1}^n X_i.$$

Note that $\hat{\mu}_1$ minimizes the squared loss

$$\sum_{i=1}^n (X_i - \mu)^2.$$

It can be shown that $\hat{\mu}_1$ is a “good” estimator if the model (1.3) holds. When (1.3) is not true, in particular when there are *outliers* (large, “wrong”, observations) (*Ausreisser*), then one has to apply a more *robust* estimator.

- The (sample) median is

$$\hat{\mu}_2 := \begin{cases} X_{((n+1)/2)} & \text{when } n \text{ odd} \\ \{X_{(n/2)} + X_{(n/2+1)}\}/2 & \text{when } n \text{ is even} \end{cases},$$

where $X_{(1)} \leq \dots \leq X_{(n)}$ are the order statistics. Note that $\hat{\mu}_2$ is a minimizer of the absolute loss

$$\sum_{i=1}^n |X_i - \mu|.$$

- The Huber estimator is

$$\hat{\mu}_3 := \arg \min_{\mu} \sum_{i=1}^n \rho(X_i - \mu), \quad (1.4)$$

where

$$\rho(x) = \begin{cases} x^2 & \text{if } |x| \leq k \\ k(2|x| - k) & \text{if } |x| > k \end{cases},$$

with $k > 0$ some given threshold.

- We finally mention the α -trimmed mean, defined, for some $0 < \alpha < 1$, as

$$\hat{\mu}_4 := \frac{1}{n - 2[n\alpha]} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} X_{(i)}.$$

Note To avoid misunderstanding, we note that e.g. in (1.4), μ is used as variable over which is minimized, whereas in (1.1), μ is a parameter. These are actually distinct concepts, but it is a general convention to abuse notation and employ the same symbol μ . When further developing the theory (see Chapter 6) we shall often introduce a new symbol for the variable, e.g., (1.4) is written as

$$\hat{\mu}_3 := \arg \min_c \sum_{i=1}^n \rho(X_i - c).$$

An example of a nonparametric estimator is the empirical distribution function

$$\hat{F}_n(\cdot) := \frac{1}{n} \#\{X_i \leq \cdot, 1 \leq i \leq n\}.$$

This is an estimator of the theoretical distribution function

$$F(\cdot) := P(X \leq \cdot).$$

Any reasonable estimator is constructed according the so-called a *plug-in principle* (*Einsetzprinzip*). That is, the parameter of interest γ is written as $\gamma = Q(F)$, with Q some given map. The empirical distribution \hat{F}_n is then “plugged in”, to obtain the estimator $T := Q(\hat{F}_n)$. (We note however that problems can arise, e.g. $Q(\hat{F}_n)$ may not be well-defined ...).

Examples are the above estimators $\hat{\mu}_1, \dots, \hat{\mu}_4$ of the location parameter μ . We define the maps

$$Q_1(F) := \int x dF(x)$$

(the mean, or point of gravity, of F), and

$$Q_2(F) := F^{-1}(1/2)$$

(the median of F), and

$$Q_3(F) := \arg \min_{\mu} \int \rho(\cdot - \mu) dF,$$

and finally

$$Q_4(F) := \frac{1}{1-2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x).$$

Then $\hat{\mu}_k$ corresponds to $Q_k(\hat{F}_n)$, $k = 1, \dots, 4$. If the model (1.2) is correct, $\hat{\mu}_1, \dots, \hat{\mu}_4$ are all estimators of μ . If the model is incorrect, each $Q_k(\hat{F}_n)$ is still an estimator of $Q_k(F)$ (assuming the latter exists), but the $Q_k(F)$ may all be different aspects of F .

1.3 Comparison of estimators: risk functions

A risk function $R(\cdot, \cdot)$ measures the loss due to the error of an estimator. The risk depends on the unknown distribution, e.g. in the location model, on μ and/or F_0 . Examples are

$$R(\mu, F_0, \hat{\mu}) := \begin{cases} \mathbf{E}_{\mu, F_0} |\hat{\mu} - \mu|^p \\ \mathbf{P}_{\mu, F_0} (|\hat{\mu} - \mu| > a) \\ \dots \end{cases}.$$

Here $p \geq 1$ and $a > 0$ are chosen by the researcher.

If $\hat{\mu}$ is an *equivariant* estimator, the above risks no longer depend on μ . An estimator $\hat{\mu} := \hat{\mu}(X_1, \dots, X_n)$ is called equivariant if

$$\hat{\mu}(X_1 + c, \dots, X_n + c) = \hat{\mu}(X_1, \dots, X_n) + c, \quad \forall c.$$

Then, writing

$$\mathbf{P}_{F_0} := \mathbf{P}_{0, F_0},$$

(and likewise for the expectation \mathbf{E}_{F_0}), we have for all $t > 0$

$$\mathbf{P}_{\mu, F_0}(\hat{\mu} - \mu \leq t) = \mathbf{P}_{F_0}(\hat{\mu} \leq t),$$

that is, the distribution of $\hat{\mu} - \mu$ does not depend on μ . We then write

$$R(\mu, F_0, \hat{\mu}) := R(F_0, \hat{\mu}) := \begin{cases} \mathbf{E}_{F_0} |\hat{\mu}|^p \\ \mathbf{P}_{F_0} (|\hat{\mu}| > a) \\ \dots \end{cases}.$$

1.4 Comparison of estimators: sensitivity

We can compare estimators with respect to their sensitivity to large errors in the data. Suppose the estimator $\hat{\mu} = \hat{\mu}_n$ is defined for each n , and is symmetric in X_1, \dots, X_n .

Influence of a single additional observation

The influence function is

$$l(x) := \hat{\mu}_{n+1}(X_1, \dots, X_n, x) - \hat{\mu}_n(X_1, \dots, X_n), \quad x \in \mathbb{R}.$$

Break down point

Let for $m \leq n$,

$$\epsilon(m) := \sup_{x_1^*, \dots, x_m^*} |\hat{\mu}(x_1^*, \dots, x_m^*, X_{m+1}, \dots, X_n)|.$$

If $\epsilon(m) := \infty$, we say that with m outliers the estimator can break down. The break down point is defined as

$$\epsilon^* := \min\{m : \epsilon(m) = \infty\}/n.$$

1.5 Confidence intervals

Consider the location model (Example 1.1.2).

Definition A subset $I = I(\mathbf{X}) \subset \mathbb{R}$, depending (only) on the data $\mathbf{X} = (X_1, \dots, X_n)$, is called a confidence set (Vertrauensbereich) for μ , at level $1 - \alpha$, if

$$\mathbb{P}_{\mu, F_0}(\mu \in I) \geq 1 - \alpha, \quad \forall \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0.$$

A confidence interval is of the form

$$I := [\underline{\mu}, \bar{\mu}],$$

where the boundaries $\underline{\mu} = \underline{\mu}(\mathbf{X})$ and $\bar{\mu} = \bar{\mu}(\mathbf{X})$ depend (only) on the data \mathbf{X} .

1.5.1 Equivalence confidence sets and tests

Let for each $\mu_0 \in \mathbb{R}$, $\phi(\mathbf{X}, \mu_0) \in \{0, 1\}$ be a test at level α for the hypothesis

$$H_{\mu_0} : \mu = \mu_0.$$

Thus, we reject H_{μ_0} if and only if $\phi(\mathbf{X}, \mu_0) = 1$, and

$$\mathbb{P}_{\mu_0, F_0}(\phi(\mathbf{X}, \mu_0) = 1) \leq \alpha.$$

Then

$$I(\mathbf{X}) := \{\mu : \phi(\mathbf{X}, \mu) = \mathbf{0}\}$$

is a $(1 - \alpha)$ -confidence set for μ .

Conversely, if $I(\mathbf{X})$ is a $(1 - \alpha)$ -confidence set for μ , then, for all μ_0 , the test $\phi(\mathbf{X}, \mu_0)$ defined as

$$\phi(\mathbf{X}, \mu_0) = \begin{cases} 1 & \text{if } \mu_0 \notin I(\mathbf{X}) \\ 0 & \text{else} \end{cases}$$

is a test at level α of H_{μ_0} .

1.6 Intermezzo: quantile functions

Let F be a distribution function. Then F is *cadlag* (continue à droite, limite à gauche). Define the quantile functions

$$q_+^F(u) := \sup\{x : F(x) \leq u\},$$

and

$$q_-^F(u) := \inf\{x : F(x) \geq u\} := F^{-1}(u).$$

It holds that

$$F(q_-^F(u)) \geq u$$

and, for all $h > 0$,

$$F(q_+^F(u) - h) \leq u.$$

Hence

$$F(q_+^F(u)-) := \lim_{h \downarrow 0} F(q_+^F(u) - h) \leq u.$$

1.7 How to construct tests and confidence sets

Consider a model class

$$\mathcal{P} := \{P_\theta : \theta \in \Theta\}.$$

Moreover, consider a space Γ , and a map

$$g : \Theta \rightarrow \Gamma, \quad g(\theta) := \gamma.$$

We think of γ as the parameter of interest (as in the plug-in principle, with $\gamma = Q(P_\theta) = g(\theta)$).

For instance, in Example 1.1.2, the parameter space is $\Theta := \{\theta = (\mu, F_0), \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0\}$, and, when μ is the parameter of interest, $g(\mu, F_0) = \mu$.

To test

$$H_{\gamma_0} : \gamma = \gamma_0,$$

we look for a *pivot* (*Tür-Angel*). This is a function $Z(\mathbf{X}, \gamma)$ depending on the data \mathbf{X} and on the parameter γ , such that for all $\theta \in \Theta$, the distribution

$$\mathbb{P}_\theta(Z(\mathbf{X}, g(\theta)) \leq \cdot) := G(\cdot)$$

does not depend on θ . We note that to find a pivot is unfortunately not always possible. However, if we *do* have a pivot $Z(\mathbf{X}, \gamma)$ with distribution G , we can compute its quantile functions

$$q_L := q_+^G\left(\frac{\alpha}{2}\right), \quad q_R := q_-^G\left(1 - \frac{\alpha}{2}\right).$$

and the test

$$\phi(\mathbf{X}, \gamma_0) := \begin{cases} 1 & \text{if } Z(\mathbf{X}, \gamma_0) \notin [q_L, q_R] \\ 0 & \text{else} \end{cases}.$$

Then the test has level α for testing H_{γ_0} , with $\gamma_0 = g(\theta_0)$:

$$\begin{aligned} \mathbb{P}_{\theta_0}(\phi(\mathbf{X}, g(\theta_0)) = 1) &= P_{\theta_0}(Z(\mathbf{X}, g(\theta_0)) > q_R) + \mathbb{P}_{\theta_0}(Z(\mathbf{X}), g(\theta_0)) < q_L) \\ &= 1 - G(q_R) + G(q_L) \leq 1 - \left(1 - \frac{\alpha}{2}\right) + \frac{\alpha}{2} = \alpha. \end{aligned}$$

As example, consider again the location model (Example 1.1.2). Let

$$\Theta := \{\theta = (\mu, F_0), \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0\},$$

with \mathcal{F}_0 a subset of the collection of symmetric distributions (see (1.2)). Let $\hat{\mu}$ be an equivariant estimator, so that the distribution of $\hat{\mu} - \mu$ does not depend on μ .

- If $\mathcal{F}_0 := \{F_0\}$ is a single distribution (i.e., the distribution F_0 is known), we take $Z(\mathbf{X}, \mu) := \hat{\mu} - \mu$ as pivot. By the equivariance, this pivot has distribution G depending only on F_0 .
- If $\mathcal{F}_0 := \{F_0(\cdot) = \Phi(\cdot/\sigma) : \sigma > 0\}$, we choose $\hat{\mu} := \bar{X}_n$ where $\bar{X}_n = \sum_{i=1}^n X_i/n$ is the sample mean. As pivot, we take

$$Z(\mathbf{X}, \mu) := \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n},$$

where $S_n^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n - 1)$ is the sample variance. Then G is the Student distribution with $n - 1$ degrees of freedom.

- If $\mathcal{F}_0 := \{F_0 \text{ continuous at } x = 0\}$, we let the pivot be the sign test statistic:

$$Z(\mathbf{X}, \mu) := \sum_{i=1}^n \mathbb{1}\{X_i \geq \mu\}.$$

Then G is the Binomial(n, p) distribution, with parameter $p = 1/2$.

Let $Z_n(X_1, \dots, X_n, \gamma)$ be some function of the data and the parameter of interest, defined for each sample size n . We call $Z_n(X_1, \dots, X_n, \gamma)$ an *asymptotic* pivot if for all $\theta \in \Theta$,

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta}(Z_n(X_1, \dots, X_n, \gamma) \leq \cdot) = G(\cdot),$$

where the limit G does not depend on θ .

In the location model, suppose X_1, \dots, X_n are the first n of an infinite sequence of i.i.d. random variables, and that

$$\mathcal{F}_0 := \{F_0 : \int x dF_0(x) = 0, \int x^2 dF_0(x) < \infty\}.$$

Then

$$Z_n(X_1, \dots, X_n, \mu) := \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n}$$

is an asymptotic pivot, with limiting distribution $G = \Phi$.

Comparison of confidence intervals and tests

When comparing confidence intervals, the aim is usually to take the one with smallest length on average (keeping the level at $1 - \alpha$). In the case of tests, we look for the one with maximal power. In the location model, this leads to studying

$$\mathbf{E}_{\mu, F_0} |\bar{\mu}(\mathbf{X}) - \underline{\mu}(\mathbf{X})|$$

for $(1 - \alpha)$ -confidence sets $[\underline{\mu}, \bar{\mu}]$, or to studying the power of test $\phi(\mathbf{X}, \mu_0)$ at level α . Recall that the power is $P_{\mu, F_0}(\phi(\mathbf{X}, \mu_0) = 1)$ for values $\mu \neq \mu_0$.

1.8 An illustration: the two-sample problem

Consider the following data, concerning weight gain/loss. The control group x had their usual diet, and the treatment group y obtained a special diet, designed for preventing weight gain. The study was carried out to test whether the diet works.

control group x	treatment group y	rank(x)	rank(y)
5	6	7	8
0	-5	3	2
16	-6	10	1
2	1	5	4
9	4	9	6
— +	— +		
32	0		

Table 2

Let n (m) be the sample size of the control group x (treatment group y). The mean in group x (y) is denoted by \bar{x} (\bar{y}). The sums of squares are $SS_x := \sum_{i=1}^n (x_i - \bar{x})^2$ and $SS_y := \sum_{j=1}^m (y_j - \bar{y})^2$. So in this study, one has $n = m = 5$ and the values $\bar{x} = 6.4$, $\bar{y} = 0$, $SS_x = 161.2$ and $SS_y = 114$. The ranks, $\text{rank}(x)$ and $\text{rank}(y)$, are the rank-numbers when putting all $n + m$ data together (e.g., $y_3 = -6$ is the smallest observation and hence $\text{rank}(y_3) = 1$).

We assume that the data are realizations of two independent samples, say $\mathbf{X} = (X_1, \dots, X_n)$ and $\mathbf{Y} = (Y_1, \dots, Y_m)$, where X_1, \dots, X_n are i.i.d. with distribution function F_X , and Y_1, \dots, Y_m are i.i.d. with distribution function F_Y . The distribution functions F_X and F_Y may be in whole or in part unknown. The testing problem is:

$$H_0 : F_X = F_Y$$

against a one- or two-sided alternative.

1.8.1 Assuming normality

The classical two-sample student test is based on the assumption that the data come from a normal distribution. Moreover, it is assumed that the variance of F_X and F_Y are equal. Thus,

$$(F_X, F_Y) \in \left\{ F_X = \Phi\left(\frac{\cdot - \mu}{\sigma}\right), F_Y = \Phi\left(\frac{\cdot - (\mu + \gamma)}{\sigma}\right) : \mu \in \mathbb{R}, \sigma > 0, \gamma \in \Gamma \right\}.$$

Here, $\Gamma \supset \{0\}$ is the range of shifts in mean one considers, e.g. $\Gamma = \mathbb{R}$ for two-sided situations, and $\Gamma = (-\infty, 0]$ for a one-sided situation. The testing problem reduces to

$$H_0 : \gamma = 0.$$

We now look for a pivot $Z(\mathbf{X}, \mathbf{Y}, \gamma)$. Define the sample means

$$\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i, \quad \bar{Y} := \frac{1}{m} \sum_{j=1}^m Y_j,$$

and the pooled sample variance

$$S^2 := \frac{1}{m+n-2} \left\{ \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (Y_j - \bar{Y})^2 \right\}.$$

Note that \bar{X} has expectation μ and variance σ^2/n , and \bar{Y} has expectation $\mu + \gamma$ and variance σ^2/m . So $\bar{Y} - \bar{X}$ has expectation γ and variance

$$\frac{\sigma^2}{n} + \frac{\sigma^2}{m} = \sigma^2 \left(\frac{n+m}{nm} \right).$$

The normality assumption implies that

$$\bar{Y} - \bar{X} \text{ is } \mathcal{N}\left(\gamma, \sigma^2 \left(\frac{n+m}{nm} \right)\right)\text{-distributed.}$$

Hence

$$\sqrt{\frac{nm}{n+m}} \left(\frac{\bar{Y} - \bar{X} - \gamma}{\sigma} \right) \text{ is } \mathcal{N}(0, 1)\text{-distributed.}$$

To arrive at a pivot, we now plug in the estimate S for the unknown σ :

$$Z(\mathbf{X}, \mathbf{Y}, \gamma) := \sqrt{\frac{nm}{n+m}} \left(\frac{\bar{Y} - \bar{X} - \gamma}{S} \right).$$

Indeed, $Z(\mathbf{X}, \mathbf{Y}, \gamma)$ has a distribution G which does not depend on unknown parameters. The distribution G is Student($n+m-2$) (the Student-distribution with $n+m-2$ degrees of freedom). As test statistic for $H_0 : \gamma = 0$, we therefore take

$$T = T^{\text{Student}} := Z(\mathbf{X}, \mathbf{Y}, 0).$$

The one-sided test at level α , for $H_0 : \gamma = 0$ against $H_1 : \gamma < 0$, is

$$\phi(\mathbf{X}, \mathbf{Y}) := \begin{cases} 1 & \text{if } T < -t_{n+m-2}(1-\alpha) \\ 0 & \text{if } T \geq -t_{n+m-2}(1-\alpha) \end{cases},$$

where, for $\nu > 0$, $t_\nu(1-\alpha) = -t_\nu(\alpha)$ is the $(1-\alpha)$ -quantile of the Student(ν)-distribution.

Let us apply this test to the data given in Table 2. We take $\alpha = 0.05$. The observed values are $\bar{x} = 6.4$, $\bar{y} = 0$ and $s^2 = 34.4$. The test statistic takes the value -1.725 which is bigger than the 5% quantile $t_8(0.05) = -1.9$. Hence, we cannot reject H_0 . The p -value of the observed value of T is

$$p\text{-value} := \mathbb{P}_{\gamma=0}(T < -1.725) = 0.06.$$

So the p -value is in this case only a little larger than the level $\alpha = 0.05$.

1.8.2 A nonparametric test

In this subsection, we suppose that F_X and F_Y are continuous, but otherwise unknown. The model class for both F_X and F_Y is thus

$$\mathcal{F} := \{\text{all continuous distributions}\}.$$

The continuity assumption ensures that all observations are distinct, that is, there are no ties. We can then put them in strictly increasing order. Let $N = n + m$ and Z_1, \dots, Z_N be the pooled sample

$$Z_i := X_i, \quad i = 1, \dots, n, \quad Z_{n+j} := Y_j, \quad j = 1, \dots, m.$$

Define

$$R_i := \text{rank}(Z_i), \quad i = 1, \dots, N.$$

and let

$$Z_{(1)} < \dots < Z_{(N)}$$

be the order statistics of the pooled sample (so that $Z_i = Z_{(R_i)}$ ($i = 1, \dots, n$)). The Wilcoxon test statistic is

$$T = T^{\text{Wilcoxon}} := \sum_{i=1}^n R_i.$$

One may check that this test statistic T can alternatively be written as

$$T = \#\{Y_j < X_i\} + \frac{n(n+1)}{2}.$$

For example, for the data in Table 2, the observed value of T is 34, and

$$\#\{y_j < x_i\} = 19, \quad \frac{n(n+1)}{2} = 15.$$

Large values of T mean that the X_i are generally larger than the Y_j , and hence indicate evidence against H_0 .

To check whether or not the observed value of the test statistic is compatible with the null-hypothesis, we need to know its null-distribution, that is, the distribution under H_0 . Under $H_0 : F_X = F_Y$, the vector of ranks (R_1, \dots, R_n) has the same distribution as n random draws without replacement from the numbers $\{1, \dots, N\}$. That is, if we let

$$\mathbf{r} := (r_1, \dots, r_n, r_{n+1}, \dots, r_N)$$

denote a permutation of $\{1, \dots, N\}$, then

$$\mathbb{P}_{H_0} \left((R_1, \dots, R_n, R_{n+1}, \dots, R_N) = \mathbf{r} \right) = \frac{1}{N!},$$

(see Theorem 1.8.1), and hence

$$\mathbb{P}_{H_0}(T = t) = \frac{\#\{\mathbf{r} : \sum_{i=1}^n r_i = t\}}{N!}.$$

This can also be written as

$$\mathbb{P}_{H_0}(T = t) = \frac{1}{\binom{N}{n}} \#\{r_1 < \dots < r_n < r_{n+1} < \dots < r_N : \sum_{i=1}^n r_i = t\}.$$

So clearly, the null-distribution of T does not depend on F_X or F_Y . It does however depend on the sample sizes n and m . It is tabulated for n and m small or moderately large. For large n and m , a normal approximation of the null-distribution can be used.

Theorem 1.8.1 formally derives the null-distribution of the test, and actually proves that the order statistics and the ranks are independent. The latter result will be of interest in Example 2.10.4.

For two random variables X and Y , use the notation

$$X \stackrel{D}{=} Y$$

when X and Y have the same distribution.

Theorem 1.8.1 *Let Z_1, \dots, Z_N be i.i.d. with continuous distribution F on \mathbb{R} . Then $(Z_{(1)}, \dots, Z_{(N)})$ and $\mathbf{R} := (R_1, \dots, R_N)$ are independent, and for all permutations $\mathbf{r} := (r_1, \dots, r_N)$,*

$$\mathbb{P}(\mathbf{R} = \mathbf{r}) = \frac{1}{N!}.$$

Proof. Let $Z_{Q_i} := Z_{(i)}$, and $\mathbf{Q} := (Q_1, \dots, Q_N)$. Then

$$\mathbf{R} = \mathbf{r} \Leftrightarrow \mathbf{Q} = \mathbf{r}^{-1} := \mathbf{q},$$

where \mathbf{r}^{-1} is the inverse permutation of \mathbf{r} .¹ For all permutations \mathbf{q} and all measurable maps f ,

$$f(Z_1, \dots, Z_N) \stackrel{\text{D}}{=} f(Z_{q_1}, \dots, Z_{q_N}).$$

Therefore, for all measurable sets $A \subset \mathbb{R}^N$, and all permutations \mathbf{q} ,

$$\begin{aligned} & \mathbb{P}\left((Z_1, \dots, Z_N) \in A, Z_1 < \dots < Z_N\right) \\ &= \mathbb{P}\left((Z_{q_1}, \dots, Z_{q_N}) \in A, Z_{q_1} < \dots < Z_{q_N}\right). \end{aligned}$$

Because there are $N!$ permutations, we see that for any \mathbf{q} ,

$$\begin{aligned} \mathbb{P}\left((Z_{(1)}, \dots, Z_{(n)}) \in A\right) &= N! \mathbb{P}\left((Z_{q_1}, \dots, Z_{q_N}) \in A, Z_{q_1} < \dots < Z_{q_N}\right) \\ &= N! \mathbb{P}\left((Z_{(1)}, \dots, Z_{(N)}) \in A, \mathbf{R} = \mathbf{r}\right), \end{aligned}$$

where $\mathbf{r} = \mathbf{q}^{-1}$. Thus we have shown that for all measurable A , and for all \mathbf{r} ,

$$\mathbb{P}\left((Z_{(1)}, \dots, Z_{(N)}) \in A, \mathbf{R} = \mathbf{r}\right) = \frac{1}{N!} \mathbb{P}\left((Z_{(1)}, \dots, Z_{(n)}) \in A\right). \quad (1.5)$$

Take $A = \mathbb{R}^N$ to find that (1.5) implies

$$\mathbb{P}\left(\mathbf{R} = \mathbf{r}\right) = \frac{1}{N!}.$$

Plug this back into (1.5) to see that we have the product structure

$$\mathbb{P}\left((Z_{(1)}, \dots, Z_{(N)}) \in A, \mathbf{R} = \mathbf{r}\right) = \mathbb{P}\left((Z_{(1)}, \dots, Z_{(n)}) \in A\right) \mathbb{P}\left(\mathbf{R} = \mathbf{r}\right),$$

which holds for all measurable A . In other words, $(Z_{(1)}, \dots, Z_{(N)})$ and \mathbf{R} are independent. \square

1.8.3 Comparison of Student's test and Wilcoxon's test

Because Wilcoxon's test is only based on the ranks, and does not rely on the assumption of normality, it lies at hand that, when the data are in fact normally distributed, Wilcoxon's test will have less power than Student's test. The loss

¹Here is an example, with $N = 3$:

$$(z_1, z_2, z_3) = (5, 6, 4)$$

$$(r_1, r_2, r_3) = (2, 3, 1)$$

$$(q_1, q_2, q_3) = (3, 1, 2)$$

of power is however small. Let us formulate this more precisely, in terms of the relative efficiency of the two tests. Let the significance α be fixed, and let β be the required power. Let n and m be equal, $N = 2n$ be the total sample size, and N^{Student} (N^{Wilcoxon}) be the number of observations needed to reach power β using Student's (Wilcoxon's) test. Consider shift alternatives, i.e. $F_Y(\cdot) = F_X(\cdot - \gamma)$, (with, in our example, $\gamma < 0$). One can show that $N^{\text{Student}}/N^{\text{Wilcoxon}}$ is approximately .95 when the normal model is correct. For a large class of distributions, the ratio $N^{\text{Student}}/N^{\text{Wilcoxon}}$ ranges from .85 to ∞ , that is, when using Wilcoxon one generally has very limited loss of efficiency as compared to Student, and one may in fact have a substantial gain of efficiency.

1.9 How to construct estimators

Consider i.i.d. observations X_1, \dots, X_n , copies of a random variable X with distribution $P \in \{P_\theta : \theta \in \Theta\}$. The parameter of interest is denoted by $\gamma = g(\theta) \in \Gamma$.

1.9.1 Plug-in estimators

For real-valued observations, one can define the distribution function

$$F(\cdot) = P(X \leq \cdot).$$

An estimator of F is the empirical distribution function

$$\hat{F}_n(\cdot) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq \cdot\}.$$

Note that when knowing only \hat{F}_n , one can reconstruct the order statistics $X_{(1)} \leq \dots \leq X_{(n)}$, but not the original data X_1, \dots, X_n . Now, the order at which the data are given carries no information about the distribution P . In other words, a “reasonable”² estimator $T = T(X_1, \dots, X_n)$ depends only on the sample (X_1, \dots, X_n) via the order statistics $(X_{(1)}, \dots, X_{(n)})$ (i.e., shuffling the data should have no influence on the value of T). Because these order statistics can be determined from the empirical distribution \hat{F}_n , we conclude that any “reasonable” estimator T can be written as a function of \hat{F}_n :

$$T = Q(\hat{F}_n),$$

for some map Q .

Similarly, the distribution function $F_\theta := P_\theta(X \leq \cdot)$ completely characterizes the distribution P . Hence, a parameter is a function of F_θ :

$$\gamma = g(\theta) = Q(F_\theta).$$

²What is “reasonable” has to be considered with some care. There are in fact “reasonable” statistical procedures that do treat the $\{X_i\}$ in an asymmetric way. An example is splitting the sample into a training and test set (for model validation).

If the mapping Q is defined at all F_θ as well as at \hat{F}_n , we call $Q(\hat{F}_n)$ a plug-in estimator of $Q(F_\theta)$.

The idea is not restricted to the one-dimensional setting. For arbitrary observation space \mathcal{X} , we define the empirical measure

$$\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i},$$

where δ_x is a point-mass at x . The empirical measure puts mass $1/n$ at each observation. This is indeed an extension of $\mathcal{X} = \mathbb{R}$ to general \mathcal{X} , as the empirical distribution function \hat{F}_n jumps at each observation, with jump height equal to the number of times the value was observed (i.e. jump height $1/n$ if all X_i are distinct). So, as in the real-valued case, if the map Q is defined at all P_θ as well as at \hat{P}_n , we call $Q(\hat{P}_n)$ a plug-in estimator of $Q(P_\theta)$.

We stress that typically, the representation $\gamma = g(\theta)$ as function Q of P_θ is not unique, i.e., that there are various choices of Q . Each such choice generally leads to a different estimator. Moreover, the assumption that Q is defined at \hat{P}_n is often violated. One can sometimes modify the map Q to a map Q_n that, in some sense, approximates Q for n large. The modified plug-in estimator then takes the form $Q_n(\hat{P}_n)$.

1.9.2 The method of moments

Let $X \in \mathbb{R}$ and suppose (say) that the parameter of interest is θ itself, and that $\Theta \subset \mathbb{R}^p$. Let $\mu_1(\theta), \dots, \mu_p(\theta)$ denote the first p moments of X (assumed to exist), i.e.,

$$\mu_j(\theta) = E_\theta X^j = \int x^j dF_\theta(x), \quad j = 1, \dots, p.$$

Also assume that the map

$$m : \Theta \rightarrow \mathbb{R}^p,$$

defined by

$$m(\theta) = [\mu_1(\theta), \dots, \mu_p(\theta)],$$

has an inverse

$$m^{-1}(\mu_1, \dots, \mu_p),$$

for all $[\mu_1, \dots, \mu_p] \in \mathcal{M}$ (say). We estimate the μ_j by their sample counterparts

$$\hat{\mu}_j := \frac{1}{n} \sum_{i=1}^n X_i^j = \int x^j d\hat{F}_n(x), \quad j = 1, \dots, p.$$

When $[\hat{\mu}_1, \dots, \hat{\mu}_p] \in \mathcal{M}$ we can plug them in to obtain the estimator

$$\hat{\theta} := m^{-1}(\hat{\mu}_1, \dots, \hat{\mu}_p).$$

Example

Let X have the negative binomial distribution with known parameter k and unknown success parameter $\theta \in (0, 1)$:

$$P_\theta(X = x) = \binom{k+x-1}{x} \theta^k (1-\theta)^x, x \in \{0, 1, \dots\}.$$

This is the distribution of the number of failures till the k^{th} success, where at each trial, the probability of success is θ , and where the trials are independent. It holds that

$$E_\theta(X) = k \frac{(1-\theta)}{\theta} := m(\theta).$$

Hence

$$m^{-1}(\mu) = \frac{k}{\mu + k},$$

and the method of moments estimator is

$$\hat{\theta} = \frac{k}{\bar{X} + k} = \frac{nk}{\sum_{i=1}^n X_i + nk} = \frac{\text{number of successes}}{\text{number of trials}}.$$

Example

Suppose X has density

$$p_\theta(x) = \theta(1+x)^{-(1+\theta)}, x > 0,$$

with respect to Lebesgue measure, and with $\theta \in \Theta \subset (0, \infty)$. Then, for $\theta > 1$

$$E_\theta X = \frac{1}{\theta - 1} := m(\theta),$$

with inverse

$$m^{-1}(\mu) = \frac{1 + \mu}{\mu}.$$

The method of moments estimator would thus be

$$\hat{\theta} = \frac{1 + \bar{X}}{\bar{X}}.$$

However, the mean $E_\theta X$ does not exist for $\theta < 1$, so when Θ contains values $\theta < 1$, the method of moments is perhaps not a good idea. We will see that the maximum likelihood estimator does not suffer from this problem.

1.9.3 Likelihood methods

Suppose that $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ is dominated by a σ -finite measure ν . We write the densities as

$$p_\theta := \frac{dP_\theta}{d\nu}, \theta \in \Theta.$$

Definition The likelihood function (of the data $\mathbf{X} = (X_1, \dots, X_n)$) is

$$L_{\mathbf{X}}(\vartheta) := \prod_{i=1}^n p_\vartheta(X_i).$$

The MLE (*maximum likelihood estimator*) is

$$\hat{\theta} := \arg \max_{\vartheta \in \Theta} L_{\mathbf{X}}(\vartheta).$$

Note We use the symbol ϑ for the variable in the likelihood function, and the slightly different symbol θ for the parameter we want to estimate. It is however a common convention to use the same symbol for both (as already noted in the earlier section on estimation). However, as we will see below, different symbols are needed for the development of the theory.

Note Alternatively, we may write the MLE as the maximizer of the *log*-likelihood

$$\hat{\theta} = \arg \max_{\vartheta \in \Theta} \log L_{\mathbf{X}}(\vartheta) = \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_{\vartheta}(X_i).$$

The log-likelihood is generally mathematically more tractable. For example, if the densities are differentiable, one can typically obtain the maximum by setting the derivatives to zero, and it is easier to differentiate a sum than a product.

Note The likelihood function may have local maxima. Moreover, the MLE is not always unique, or may not exist (for example, the likelihood function may be unbounded).

We will now show that maximum likelihood is a plug-in method. First, as noted above, the MLE maximizes the log-likelihood. We may of course normalize the log-likelihood by $1/n$:

$$\hat{\theta} = \arg \max_{\vartheta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p_{\vartheta}(X_i).$$

Replacing the average $\sum_{i=1}^n \log p_{\vartheta}(X_i)/n$ by its theoretical counterpart gives

$$\arg \max_{\vartheta \in \Theta} E_{\theta} \log p_{\vartheta}(X)$$

which is indeed equal to the parameter θ we are trying to estimate: by the inequality $\log x \leq x - 1$, $x > 0$,

$$E_{\theta} \log \frac{p_{\vartheta}(X)}{p_{\theta}(X)} \leq E_{\theta} \left(\frac{p_{\vartheta}(X)}{p_{\theta}(X)} - 1 \right) = 0.$$

(Note that using different symbols ϑ and θ is indeed crucial here.) Chapter 6 will put this in a wider perspective.

Example

We turn back to the previous example. Suppose X has density

$$p_{\theta}(x) = \theta(1+x)^{-(1+\theta)}, \quad x > 0,$$

with respect to Lebesgue measure, and with $\theta \in \Theta = (0, \infty)$. Then

$$\begin{aligned}\log p_\vartheta(x) &= \log \vartheta - (1 + \vartheta) \log(1 + x), \\ \frac{d}{d\vartheta} \log p_\vartheta(x) &= \frac{1}{\vartheta} - \log(1 + x).\end{aligned}$$

We put the derivative of the log-likelihood to zero and solve:

$$\begin{aligned}\frac{n}{\hat{\theta}} - \sum_{i=1}^n \log(1 + X_i) &= 0 \\ \Rightarrow \hat{\theta} &= \frac{1}{\{\sum_{i=1}^n \log(1 + X_i)\}/n}.\end{aligned}$$

(One may check that this is indeed the maximum.)

Example

Let $X \in \mathbb{R}$ and $\theta = (\mu, \sigma^2)$, with $\mu \in \mathbb{R}$ a location parameter, $\sigma > 0$ a scale parameter. We assume that the distribution function F_θ of X is

$$F_\theta(\cdot) = F_0\left(\frac{\cdot - \mu}{\sigma}\right),$$

where F_0 is a given distribution function, with density f_0 w.r.t. Lebesgue measure. The density of X is thus

$$p_\theta(\cdot) = \frac{1}{\sigma} f_0\left(\frac{\cdot - \mu}{\sigma}\right).$$

Case 1 If $F_0 = \Phi$ (the standard normal distribution), then

$$f_0(x) = \phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}x^2\right], \quad x \in \mathbb{R},$$

so that

$$p_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right], \quad x \in \mathbb{R}.$$

The MLE of μ resp. σ^2 is

$$\hat{\mu} = \bar{X}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Case 2 The (standardized) *double exponential* or *Laplace* distribution has density

$$f_0(x) = \frac{1}{\sqrt{2}} \exp\left[-\sqrt{2}|x|\right], \quad x \in \mathbb{R},$$

so

$$p_\theta(x) = \frac{1}{\sqrt{2}\sigma^2} \exp\left[-\frac{\sqrt{2}|x - \mu|}{\sigma}\right], \quad x \in \mathbb{R}.$$

The MLE of μ resp. σ is now

$$\hat{\mu} = \text{sample median}, \quad \hat{\sigma} = \frac{\sqrt{2}}{n} \sum_{i=1}^n |X_i - \hat{\mu}_2|.$$

Example

Here is a famous example, from Kiefer and Wolfowitz (1956), where the likelihood is unbounded, and hence the MLE does not exist. It concerns the case of a mixture of two normals: each observation, is either $\mathcal{N}(\mu, 1)$ -distributed or $\mathcal{N}(\mu, \sigma^2)$ -distributed, each with probability 1/2 (say). The unknown parameter is $\theta = (\mu, \sigma^2)$, and X has density

$$p_{\theta}(x) = \frac{1}{2}\phi(x - \mu) + \frac{1}{2\sigma}\phi((x - \mu)/\sigma), \quad x \in \mathbb{R},$$

w.r.t. Lebesgue measure. Then

$$L_{\mathbf{X}}(\mu, \sigma^2) = \prod_{i=1}^n \left(\frac{1}{2}\phi(X_i - \mu) + \frac{1}{2\sigma}\phi((X_i - \mu)/\sigma) \right).$$

Taking $\mu = X_1$ yields

$$L_{\mathbf{X}}(X_1, \sigma^2) = \frac{1}{\sqrt{2\pi}} \left(\frac{1}{2} + \frac{1}{2\sigma} \right) \prod_{i=2}^n \left(\frac{1}{2}\phi(X_i - X_1) + \frac{1}{2\sigma}\phi((X_i - X_1)/\sigma) \right).$$

Now, since for all $z \neq 0$

$$\lim_{\sigma \downarrow 0} \frac{1}{\sigma} \phi(z/\sigma) = 0,$$

we have

$$\lim_{\sigma \downarrow 0} \prod_{i=2}^n \left(\frac{1}{2}\phi(X_i - X_1) + \frac{1}{2\sigma}\phi((X_i - X_1)/\sigma) \right) = \prod_{i=2}^n \frac{1}{2}\phi(X_i - X_1) > 0.$$

It follows that

$$\lim_{\sigma \downarrow 0} L_{\mathbf{X}}(X_1, \sigma^2) = \infty.$$

Asymptotic tests and confidence intervals based on the likelihood

Suppose that Θ is an open subset of \mathbb{R}^p . Define the log-likelihood ratio

$$Z(\mathbf{X}, \theta) := 2 \left\{ \log L_{\mathbf{X}}(\hat{\theta}) - \log L_{\mathbf{X}}(\theta) \right\}.$$

Note that $Z(\mathbf{X}, \theta) \geq 0$, as $\hat{\theta}$ maximizes the (log)-likelihood. We will see in Chapter 6 that, under some regularity conditions,

$$Z(\mathbf{X}, \theta) \xrightarrow{D_{\theta}} \chi_p^2, \quad \forall \theta.$$

Here, “ $\xrightarrow{D_\theta}$ ” means convergence in distribution under \mathbb{P}_θ , and χ_p^2 denotes the Chi-squared distribution with p degrees of freedom.

Thus, $Z(\mathbf{X}, \theta)$ is an asymptotic pivot. For the null-hypotheses

$$H_0 : \theta = \theta_0,$$

a test at asymptotic level α is: reject H_0 if $Z(\mathbf{X}, \theta_0) > \chi_p^2(1-\alpha)$, where $\chi_p^2(1-\alpha)$ is the $(1-\alpha)$ -quantile of the χ_p^2 -distribution. An asymptotic $(1-\alpha)$ -confidence set for θ is

$$\begin{aligned} & \{\theta : Z(\mathbf{X}, \theta) \leq \chi_p^2(1-\alpha)\} \\ &= \{\theta : 2 \log L_{\mathbf{X}}(\hat{\theta}) \leq 2 \log L_{\mathbf{X}}(\theta) + \chi_p^2(1-\alpha)\}. \end{aligned}$$

Example

Here is a toy example. Let X have the $\mathcal{N}(\mu, 1)$ -distribution, with $\mu \in \mathbb{R}$ unknown. The MLE of μ is the sample average $\hat{\mu} = \bar{X}$. It holds that

$$\log L_{\mathbf{X}}(\hat{\mu}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})^2,$$

and

$$2 \left\{ \log L_{\mathbf{X}}(\hat{\mu}) - \log L_{\mathbf{X}}(\mu) \right\} = n(\bar{X} - \mu)^2.$$

The random variable $\sqrt{n}(\bar{X} - \mu)$ is $\mathcal{N}(0, 1)$ -distributed under \mathbb{P}_μ . So its square, $n(\bar{X} - \mu)^2$, has a χ_1^2 -distribution. Thus, in this case the above test (confidence interval) is exact.

...

Chapter 2

Decision theory

Notation

In this chapter, we denote the observable random variable (the data) by $X \in \mathcal{X}$, and its distribution by $P \in \mathcal{P}$. The probability model is $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$, with θ an unknown parameter. In particular cases, we apply the results with X being replaced by a vector $\mathbf{X} = (X_1, \dots, X_n)$, with X_1, \dots, X_n i.i.d. with distribution $P \in \{P_\theta : \theta \in \Theta\}$ (so that \mathbf{X} has distribution $\mathbb{P} := \prod_{i=1}^n P \in \{\mathbb{P}_\theta = \prod_{i=1}^n P_\theta : \theta \in \Theta\}$).

2.1 Decisions and their risk

Let \mathcal{A} be the *action space*.

- $\mathcal{A} = \mathbb{R}$ corresponds to estimating a real-valued parameter.
- $\mathcal{A} = \{0, 1\}$ corresponds to testing a hypothesis.
- $\mathcal{A} = [0, 1]$ corresponds to randomized tests.
- $\mathcal{A} = \{\text{intervals}\}$ corresponds to confidence intervals.

Given the observation X , we decide to take a certain action in \mathcal{A} . Thus, an action is a map $d : \mathcal{X} \rightarrow \mathcal{A}$, with $d(X)$ being the decision taken.

A *loss function* (*Verlustfunktion*) is a map

$$L : \Theta \times \mathcal{A} \rightarrow \mathbb{R},$$

with $L(\theta, a)$ being the loss when the parameter value is θ and one takes action a .

The risk of decision $d(X)$ is defined as

$$R(\theta, d) := E_\theta L(\theta, d(X)), \theta \in \Theta.$$

Example 2.1.1 In the case of estimating a parameter of interest $g(\theta) \in \mathbb{R}$, the action space is $\mathcal{A} = \mathbb{R}$ (or a subset thereof). Important loss functions are then

$$L(\theta, a) := w(\theta)|g(\theta) - a|^r,$$

where $w(\cdot)$ are given non-negative weights and $r \geq 0$ is a given power. The risk is then

$$R(\theta, d) = w(\theta)E_\theta|g(\theta) - d(X)|^r.$$

A special case is taking $w \equiv 1$ and $r = 2$. Then

$$R(\theta, d) = E_\theta|g(\theta) - d(X)|^2$$

is called the *mean square error*.

Example 2.1.2 Consider testing the hypothesis

$$H_0 : \theta \in \Theta_0$$

against the alternative

$$H_1 : \theta \in \Theta_1.$$

Here, Θ_0 and Θ_1 are given subsets of Θ with $\Theta_0 \cap \Theta_1 = \emptyset$. As action space, we take $\mathcal{A} = \{0, 1\}$, and as loss

$$L(\theta, a) := \begin{cases} 1 & \text{if } \theta \in \Theta_0 \text{ and } a = 1 \\ c & \text{if } \theta \in \Theta_1 \text{ and } a = 0 \\ 0 & \text{otherwise} \end{cases}.$$

Here $c > 0$ is some given constant. Then

$$R(\theta, d) = \begin{cases} P_\theta(d(X) = 1) & \text{if } \theta \in \Theta_0 \\ cP_\theta(d(X) = 0) & \text{if } \theta \in \Theta_1 \\ 0 & \text{otherwise} \end{cases}.$$

Thus, the risks correspond to the error probabilities (type I and type II errors).

Note

The best decision d is the one with the smallest risk $R(\theta, d)$. However, θ is not known. Thus, if we compare two decision functions d_1 and d_2 , we may run into problems because the risks are not comparable: $R(\theta, d_1)$ may be smaller than $R(\theta, d_2)$ for some values of θ , and larger than $R(\theta, d_2)$ for other values of θ .

Example 2.1.3 Let $X \in \mathbb{R}$ and let $g(\theta) = E_\theta X := \mu$. We take quadratic loss

$$L(\theta, a) := |\mu - a|^2.$$

Assume that $\text{var}_\theta(X) = 1$ for all θ . Consider the collection of decisions

$$d_\lambda(X) := \lambda X,$$

where $0 \leq \lambda \leq 1$. Then

$$R(\theta, d_\lambda) = \text{var}(\lambda X) + \text{bias}_\theta^2(\lambda X)$$

$$= \lambda^2 + (\lambda - 1)^2 \mu^2.$$

The “optimal” choice for λ would be

$$\lambda_{\text{opt}} := \frac{\mu^2}{1 + \mu^2},$$

because this value minimizes $R(\theta, d_\lambda)$. However, λ_{opt} depends on the unknown μ , so $d_{\lambda_{\text{opt}}}(X)$ is not an estimator.

Various optimality concepts

We will consider three optimality concepts: *admissibility* (*zulässigkeit*), *mini-max* and *Bayes*.

2.2 Admissibility

Definition A decision d' is called strictly better than d if

$$R(\theta, d') \leq R(\theta, d), \quad \forall \theta,$$

and

$$\exists \theta : R(\theta, d') < R(\theta, d).$$

When there exists a d' that is strictly better than d , then d is called inadmissible.

Example 2.2.1 Let, for $n \geq 2$, X_1, \dots, X_n be i.i.d., with $g(\theta) := E_\theta(X_i) := \mu$, and $\text{var}(X_i) = 1$ (for all i). Take quadratic loss $L(\theta, a) := |\mu - a|^2$. Consider $d'(X_1, \dots, X_n) := \bar{X}_n$ and $d(X_1, \dots, X_n) := X_1$. Then, $\forall \theta$,

$$R(\theta, d') = \frac{1}{n}, \quad R(\theta, d) = 1,$$

so that d is inadmissible.

Note

We note that to show that a decision d is inadmissible, it suffices to find a strictly better d' . On the other hand, to show that d is admissible, one has to verify that there is no strictly better d' . So in principle, one then has to take all possible d' into account.

Example 2.2.2 Let $L(\theta, a) := |g(\theta) - a|^r$ and $d(X) := g(\theta_0)$, where θ_0 is some fixed given value.

Lemma Assume that P_{θ_0} dominates P_θ ¹ for all θ . Then d is admissible.

Proof.

¹Let P and Q be probability measures on the same measurable space. Then P dominates Q if for all measurable B , $P(B) = 0$ implies $Q(B) = 0$ (Q is absolut stetig bezüglich P).

Suppose that d' is better than d . Then we have

$$E_{\theta_0}|g(\theta_0) - d'(X)|^r \leq 0.$$

This implies that

$$d'(X) = g(\theta_0), \quad P_{\theta_0}\text{-almost surely.} \quad (2.1)$$

Since by (2.1),

$$P_{\theta_0}(d'(X) \neq g(\theta_0)) = 0$$

the assumption that P_{θ_0} dominates P_{θ} , $\forall \theta$, implies now

$$P_{\theta}(d'(X) \neq g(\theta_0)) = 0, \quad \forall \theta.$$

That is, for all θ , $d'(X) = g(\theta_0)$, P_{θ} -almost surely, and hence, for all θ , $R(\theta, d') = R(\theta, d)$. So d' is not strictly better than d . We conclude that d is admissible. \square

Example 2.2.3 We consider testing

$$H_0 : \theta = \theta_0$$

against the alternative

$$H_1 : \theta = \theta_1.$$

We let $\mathcal{A} = [0, 1]$ and let $d := \phi$ be a randomized test. As risk, we take

$$R(\theta, \phi) := \begin{cases} E_{\theta}\phi(X), & \theta = \theta_0 \\ 1 - E_{\theta}\phi(X), & \theta = \theta_1 \end{cases}.$$

We let p_0 (p_1) be the density of P_{θ_0} (P_{θ_1}) with respect to some dominating measure ν (for example $\nu = P_{\theta_0} + P_{\theta_1}$). A Neyman Pearson test is

$$\phi_{\text{NP}} := \begin{cases} 1 & \text{if } p_1/p_0 > c \\ q & \text{if } p_1/p_0 = c \\ 0 & \text{if } p_1/p_0 < c \end{cases}.$$

Here $0 \leq q \leq 1$, and $0 \leq c < \infty$ are given constants. To check whether ϕ_{NP} is admissible, we first recall the Neyman Pearson Lemma.

Neyman Pearson Lemma *Let ϕ be some test. We have*

$$R(\theta_1, \phi_{\text{NP}}) - R(\theta_1, \phi) \leq c[R(\theta_0, \phi) - R(\theta_0, \phi_{\text{NP}})].$$

Proof.

$$\begin{aligned} R(\theta_1, \phi_{\text{NP}}) - R(\theta_1, \phi) &= \int (\phi - \phi_{\text{NP}})p_1 \\ &= \int_{p_1/p_0 > c} (\phi - \phi_{\text{NP}})p_1 + \int_{p_1/p_0 = c} (\phi - \phi_{\text{NP}})p_1 + \int_{p_1/p_0 < c} (\phi - \phi_{\text{NP}})p_1 \\ &\leq c \int_{p_1/p_0 > c} (\phi - \phi_{\text{NP}})p_0 + c \int_{p_1/p_0 = c} (\phi - \phi_{\text{NP}})p_0 + c \int_{p_1/p_0 < c} (\phi - \phi_{\text{NP}})p_0 \end{aligned}$$

$$= c[R(\theta_0, \phi) - R(\theta_0, \phi_{\text{NP}})].$$

□

Lemma *A Neyman Pearson test is admissible if and only if one of the following two cases hold:*

i) its power is strictly less than 1,

or

ii) it has minimal level among all tests with power 1.

Proof. Suppose $R(\theta_0, \phi) < R(\theta_0, \phi_{\text{NP}})$. Then from the Neyman Pearson Lemma, we now that either $R(\theta_1, \phi) > R(\theta_1, \phi_{\text{NP}})$ (i.e., then ϕ is not better than ϕ_{NP}), or $c = 0$. But when $c = 0$, it holds that $R(\theta_1, \phi_{\text{NP}}) = 0$, i.e. then ϕ_{NP} has power one.

Similarly, suppose that $R(\theta_1, \phi) < R(\theta_1, \phi_{\text{NP}})$. Then it follows from the Neyman Pearson Lemma that $R(\theta_0, \phi) > R(\theta_0, \phi_{\text{NP}})$, because we assume $c < \infty$.

□

2.3 Minimality

Definition *A decision d is called minimax if*

$$\sup_{\theta} R(\theta, d) = \inf_{d'} \sup_{\theta} R(\theta, d').$$

Thus, the minimax criterion concerns the best decision in the worst possible case.

Lemma *A Neyman Pearson test ϕ_{NP} is minimax, if and only if $R(\theta_0, \phi_{\text{NP}}) = R(\theta_1, \phi_{\text{NP}})$.*

Proof. Let ϕ be a test, and write for $j = 0, 1$,

$$r_j := R(\theta_j, \phi_{\text{NP}}), \quad r'_j = R(\theta_j, \phi).$$

Suppose that $r_0 = r_1$ and that ϕ_{NP} is not minimax. Then, for some test ϕ ,

$$\max_j r'_j < \max_j r_j.$$

This implies that both

$$r'_0 < r_0, \quad r'_1 < r_1$$

and by the Neyman Pearson Lemma, this is not possible.

Let $S = \{(R(\theta_0, \phi), R(\theta_1, \phi)) : \phi : \mathcal{X} \rightarrow [0, 1]\}$. Note that S is convex. Thus, if $r_0 < r_1$, we can find a test ϕ with $r_0 < r'_0 < r_1$ and $r'_1 < r_1$. So then ϕ_{NP} is not minimax. Similarly if $r_0 > r_1$.

□

2.4 Bayes decisions

Suppose the parameter space Θ is a measurable space. We can then equip it with a probability measure Π . We call Π the *a priori* distribution.

Definition *The Bayes risk (with respect to the probability measure Π) is*

$$r(\Pi, d) := \int_{\Theta} R(\vartheta, d) d\Pi(\vartheta).$$

A decision d is called Bayes (with respect to Π) if

$$r(\Pi, d) = \inf_{d'} r(\Pi, d').$$

If Π has density $w := d\Pi/d\mu$ with respect to some dominating measure μ , we may write

$$r(\Pi, d) = \int_{\Theta} R(\vartheta, d) w(\vartheta) d\mu(\vartheta) := r_w(d).$$

Thus, the Bayes risk may be thought of as taking a weighted average of the risks. For example, one may want to assign more weight to “important” values of θ .

Example 2.4.1 Consider again the testing problem

$$H_0 : \theta = \theta_0$$

against the alternative

$$H_1 : \theta = \theta_1.$$

Let

$$r_w(\phi) := w_0 R(\theta_0, \phi) + w_1 R(\theta_1, \phi).$$

We take $0 < w_0 = 1 - w_1 < 1$.

Lemma *Bayes test is*

$$\phi_{\text{Bayes}} = \begin{cases} 1 & \text{if } p_1/p_0 > w_0/w_1 \\ q & \text{if } p_1/p_0 = w_0/w_1 \\ 0 & \text{if } p_1/p_0 < w_0/w_1 \end{cases}.$$

Proof.

$$\begin{aligned} r_w(\phi) &= w_0 \int \phi p_0 + w_1 (1 - \int \phi p_1) \\ &= \int \phi (w_0 p_0 - w_1 p_1) + w_1. \end{aligned}$$

So we choose $\phi \in [0, 1]$ to minimize $\phi (w_0 p_0 - w_1 p_1)$. This is done by taking

$$\phi = \begin{cases} 1 & \text{if } w_0 p_0 - w_1 p_1 < 0 \\ q & \text{if } w_0 p_0 - w_1 p_1 = 0 \\ 0 & \text{if } w_0 p_0 - w_1 p_1 > 0 \end{cases},$$

where for q we may take any value between 0 and 1. □

Note that

$$2r_w(\phi_{\text{Bayes}}) = 1 - \int |w_1 p_1 - w_0 p_0|.$$

In particular, when $w_0 = w_1 = 1/2$,

$$2r_w(\phi_{\text{Bayes}}) = 1 - \int |p_1 - p_0|/2,$$

i.e., the risk is large if the two densities are close to each other.

2.5 Intermezzo: conditional distributions

Recall the definition of conditional probabilities: for two sets A and B , with $P(B) \neq 0$, the conditional probability of A given B is defined as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

It follows that

$$P(B|A) = P(A|B) \frac{P(B)}{P(A)},$$

and that, for a partition $\{B_j\}^2$

$$P(A) = \sum_j P(A|B_j)P(B_j).$$

Consider now two random vectors $X \in \mathbb{R}^n$ and $Y \in \mathbb{R}^m$. Let $f_{X,Y}(\cdot, \cdot)$, be the density of (X, Y) with respect to Lebesgue measure (assumed to exist). The marginal density of X is

$$f_X(\cdot) = \int f_{X,Y}(\cdot, y) dy,$$

and the marginal density of Y is

$$f_Y(\cdot) = \int f_{X,Y}(x, \cdot) dx.$$

Definition *The conditional density of X given $Y = y$ is*

$$f_X(x|y) := \frac{f_{X,Y}(x, y)}{f_Y(y)}, \quad x \in \mathbb{R}^n.$$

² $\{B_j\}$ is a partition if $B_j \cap B_k = \emptyset$ for all $j \neq k$ and $P(\cup_j B_j) = 1$.

Thus, we have

$$f_Y(y|x) = f_X(x|y) \frac{f_Y(y)}{f_X(x)}, \quad (x, y) \in \mathbb{R}^{n+m},$$

and

$$f_X(x) = \int f_X(x|y)f_Y(y)dy, \quad x \in \mathbb{R}^n.$$

Definition The conditional expectation of $g(X, Y)$ given $Y = y$ is

$$E[g(X, Y)|Y = y] := \int f_X(x|y)g(x, y)dx.$$

Note thus that

$$E[g_1(X)g_2(Y)|Y = y] = g_2(y)E[g_1(X)|Y = y].$$

Notation We define the random variable $E[g(X, Y)|Y]$ as

$$E[g(X, Y)|Y] := h(Y),$$

where $h(y)$ is the function $h(y) := E[g(X, Y)|Y = y]$.

Lemma 2.5.1 (Iterated expectations lemma) It holds that

$$E\left[E[g(X, Y)|Y]\right] = Eg(X, Y).$$

Proof. Define

$$h(y) := E[g(X, Y)|Y = y].$$

Then

$$\begin{aligned} Eh(Y) &= \int h(y)f_Y(y)dy = \int E[g(X, Y)|Y = y]f_Y(y)dy \\ &= \int \int g(x, y)f_{X,Y}(x, y)dxdy = Eg(X, Y). \end{aligned}$$

□

2.6 Bayes methods

Let X have distribution $P \in \mathcal{P} := \{P_\theta : \theta \in \Theta\}$. Suppose \mathcal{P} is dominated by a (σ -finite) measure ν , and let $p_\theta = dP_\theta/d\nu$ denote the densities. Let Π be an a priori distribution on Θ , with density $w := d\Pi/d\mu$. We now think of p_θ as the density of X given the value of θ . We write it as

$$p_\theta(x) = p(x|\theta), \quad x \in \mathcal{X}.$$

Moreover, we define

$$p(\cdot) := \int_{\Theta} p(\cdot|\vartheta)w(\vartheta)d\mu(\vartheta).$$

Definition The a posteriori density of θ is

$$w(\vartheta|x) = p(x|\vartheta)\frac{w(\vartheta)}{p(x)}, \quad \vartheta \in \Theta, \quad x \in \mathcal{X}.$$

Lemma 2.6.1 Given the data $X = x$, consider θ as a random variable with density $w(\vartheta|x)$. Let

$$l(x, a) := E[L(\theta, a)|X = x] := \int_{\Theta} L(\vartheta, a)w(\vartheta|x)d\mu(\vartheta),$$

and

$$d(x) := \arg \min_a l(x, a).$$

Then d is Bayes decision d_{Bayes} .

Proof.

$$\begin{aligned} r_w(d') &= \int_{\Theta} R(\vartheta, d')w(\vartheta)d\mu(\vartheta) \\ &= \int_{\Theta} \left[\int_{\mathcal{X}} L(\vartheta, d'(x))p(x|\vartheta)d\nu(x) \right] w(\vartheta)d\mu(\vartheta) \\ &= \int_{\mathcal{X}} \left[\int_{\Theta} L(\vartheta, d'(x))w(\vartheta|x)d\mu(\vartheta) \right] p(x)d\nu(x) \\ &= \int_{\mathcal{X}} l(x, d'(x))p(x)d\nu(x) \\ &\geq \int_{\mathcal{X}} l(x, d(x))p(x)d\nu(x) \\ &= r_w(d). \end{aligned}$$

□

Example 2.6.1 For the testing problem

$$H_0 : \theta = \theta_0$$

against the alternative

$$H_1 : \theta = \theta_1,$$

we have

$$l(x, \phi) = \phi w_0 p_0(x)/p(x) + (1 - \phi)w_1 p_1(x)/p(x).$$

Thus,

$$\arg \min_{\phi} l(\cdot, \phi) = \begin{cases} 1 & \text{if } w_1 p_1 > w_0 p_0 \\ \phi & \text{if } w_1 p_1 = w_0 p_0 \\ 0 & \text{if } w_1 p_1 < w_0 p_0 \end{cases}.$$

In the next example, we shall use:

Lemma. *Let Z be a real-valued random variable. Then*

$$\arg \min_{a \in \mathbb{R}} E(Z - a)^2 = EZ.$$

Proof.

$$E(Z - a)^2 = \text{var}(Z) + (a - EZ)^2.$$

□

Example 2.6.2 Consider the case $\mathcal{A} = \mathbb{R}$ and $\Theta \subseteq \mathbb{R}$. Let $L(\theta, a) := |\theta - a|^2$. Then

$$d_{\text{Bayes}}(X) = E(\theta|X).$$

Example 2.6.3 Consider again the case $\Theta \subseteq \mathbb{R}$, and $\mathcal{A} = \Theta$, and now with loss function $L(\theta, a) := 1\{|\theta - a| > c\}$ for a given constant $c > 0$. Then

$$l(x, a) = \Pi(|\theta - a| > c | X = x) = \int_{|\vartheta - a| > c} w(\vartheta|x) d\vartheta.$$

We note that for $c \rightarrow 0$

$$\frac{1 - l(x, a)}{2c} = \frac{\Pi(|\theta - a| \leq c | X = x)}{2c} \approx w(a|x) = p(x|a) \frac{w(a)}{p(x)}.$$

Thus, for c small, Bayes rule is approximately $d_0(x) := \arg \max_{a \in \Theta} p(x|a)w(a)$. The estimator $d_0(X)$ is called the maximum a posteriori estimator. If w is the uniform density on Θ (which only exists if Θ is bounded), then $d_0(X)$ is the maximum likelihood estimator.

Example 2.6.4 Suppose that given θ , X has Poisson distribution with parameter θ , and that θ has the Gamma(k, λ)-distribution. The density of θ is then

$$w(\vartheta) = \lambda^k \vartheta^{k-1} e^{-\lambda\vartheta} / \Gamma(k),$$

where

$$\Gamma(k) = \int_0^\infty e^{-z} z^{k-1} dz.$$

The Gamma(k, λ) distribution has mean

$$E\theta = \int_0^\infty \vartheta w(\vartheta) d\vartheta = \frac{k}{\lambda}.$$

The a posteriori density is then

$$\begin{aligned} w(\vartheta|x) &= p(x|\vartheta) \frac{w(\vartheta)}{p(x)} \\ &= e^{-\vartheta} \frac{\vartheta^x}{x!} \frac{\lambda^k \vartheta^{k-1} e^{-\lambda\vartheta} / \Gamma(k)}{p(x)} \end{aligned}$$

$$= e^{-\vartheta(1+\lambda)} \vartheta^{k+x-1} c(x, k, \lambda),$$

where $c(x, k, \lambda)$ is such that

$$\int w(\vartheta|x) d\vartheta = 1.$$

We recognize $w(\vartheta|x)$ as the density of the Gamma($k + x, 1 + \lambda$)-distribution. Bayes estimator with quadratic loss is thus

$$E(\theta|X) = \frac{k + X}{1 + \lambda}.$$

The maximum a posteriori estimator is

$$\frac{k + X - 1}{1 + \lambda}.$$

Example 2.6.5 Suppose given θ , X has the Binomial(n, θ)-distribution, and that θ is uniformly distributed on $[0, 1]$. Then

$$w(\vartheta|x) = \binom{n}{x} \vartheta^x (1 - \vartheta)^{n-x} / p(x).$$

This is the density of the Beta($x+1, n-x+1$)-distribution. Thus, with quadratic loss, Bayes estimator is

$$E(\theta|X) = \frac{X + 1}{n + 2}.$$

(...To be extended to general Beta prior.)

2.7 Discussion of Bayesian approach (to be written)

...

2.8 Integrating parameters out (to be written)

...

2.9 Intermezzo: some distribution theory

2.9.1 The multinomial distribution

In a survey, people were asked their opinion about some political issue. Let X be the number of **yes**-answers, Y be the number of **no** and Z be the number of **perhaps**. The total number of people in the survey is $n = X + Y + Z$. We

consider the votes as a sample with replacement, with $p_1 = P(\text{yes})$, $p_2 = P(\text{no})$, and $p_3 = P(\text{perhaps})$, $p_1 + p_2 + p_3 = 1$. Then

$$P(X = x, Y = y, Z = z) = \binom{n}{x \ y \ z} p_1^x p_2^y p_3^z, \quad (x, y, z) \in \{0, \dots, n\}, \quad x + y + z = n.$$

Here

$$\binom{n}{x \ y \ z} := \frac{n!}{x!y!z!}.$$

It is called a *multinomial* coefficient.

Lemma *The marginal distribution of X is the Binomial(n, p_1)-distribution.*

Proof. For $x \in \{0, \dots, n\}$, we have

$$\begin{aligned} P(X = x) &= \sum_{y=0}^{n-x} P(X = x, Y = y, Z = n - x - y) \\ &= \sum_{y=0}^{n-x} \binom{n}{x \ y \ n-x-y} p_1^x p_2^y (1 - p_1 - p_2)^{n-x-y} \\ &= \binom{n}{x} p_1^x \sum_{y=0}^{n-x} \binom{n-x}{y} p_2^y (1 - p_1 - p_2)^{n-x-y} = \binom{n}{x} p_1^x (1 - p_1)^{n-x}. \end{aligned}$$

□

Definition *We say that the random vector (N_1, \dots, N_k) has the multinomial distribution with parameters n and p_1, \dots, p_k (with $\sum_{j=1}^k p_j = 1$), if for all $(n_1, \dots, n_k) \in \{0, \dots, n\}^k$, with $n_1 + \dots + n_k = n$, it holds that*

$$P(N_1 = n_1, \dots, N_k = n_k) = \binom{n}{n_1 \ \dots \ n_k} p_1^{n_1} \dots p_k^{n_k}.$$

Here

$$\binom{n}{n_1 \ \dots \ n_k} := \frac{n!}{n_1! \dots n_k!}.$$

Example 2.9.1 Let X_1, \dots, X_n be i.i.d. copies of a random variable $X \in \mathbb{R}$ with distribution F , and let $-\infty = a_0 < a_1 < \dots < a_{k-1} < a_k = \infty$. Define, for $j = 1, \dots, k$,

$$p_j := P(X \in (a_{j-1}, a_j]) = F(a_j) - F(a_{j-1}),$$

$$\frac{N_j}{n} := \frac{\#\{X_i \in (a_{j-1}, a_j]\}}{n} = \hat{F}_n(a_j) - \hat{F}_n(a_{j-1}).$$

Then (N_1, \dots, N_k) has the Multinomial(n, p_1, \dots, p_k)-distribution.

2.9.2 The Poisson distribution

Definition A random variable $X \in \{0, 1, \dots\}$ has the Poisson distribution with parameter $\lambda > 0$, if for all $x \in \{0, 1, \dots\}$

$$P(X = x) = e^{-\lambda} \frac{\lambda^x}{x!}.$$

Lemma Suppose X and Y are independent, and that X has the Poisson(λ)-distribution, and Y the Poisson(μ)-distribution. Then $Z := X + Y$ has the Poisson($\lambda + \mu$)-distribution.

Proof. For all $z \in \{0, 1, \dots\}$, we have

$$\begin{aligned} P(Z = z) &= \sum_{x=0}^z P(X = x, Y = z - x) \\ &= \sum_{x=0}^z P(X = x)P(Y = z - x) = \sum_{x=0}^z e^{-\lambda} \frac{\lambda^x}{x!} e^{-\mu} \frac{\mu^{z-x}}{(z-x)!} \\ &= e^{-(\lambda+\mu)} \frac{1}{z!} \sum_{x=0}^{z-x} \binom{n}{x} \lambda^x \mu^{z-x} = e^{-(\lambda+\mu)} \frac{(\lambda + \mu)^z}{z!}. \end{aligned}$$

□

Lemma Let X_1, \dots, X_n be independent, and (for $i = 1, \dots, n$), let X_i have the Poisson(λ_i)-distribution. Define $Z := \sum_{i=1}^n X_i$. Let $z \in \{0, 1, \dots\}$. Then the conditional distribution of (X_1, \dots, X_n) given $Z = z$ is the multinomial distribution with parameters z and p_1, \dots, p_n , where

$$p_j = \frac{\lambda_j}{\sum_{i=1}^n \lambda_i}, \quad j = 1, \dots, n.$$

Proof. First note that Z is Poisson(λ_+)-distributed, with $\lambda_+ := \sum_{i=1}^n \lambda_i$. Thus, for all $(x_1, \dots, x_n) \in \{0, 1, \dots, z\}^n$ satisfying $\sum_{i=1}^n x_i = z$, we have

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n | Z = z) &= \frac{P(X_1 = x_1, \dots, X_n = x_n)}{P(Z = z)} \\ &= \frac{\prod_{i=1}^n (e^{-\lambda_i} \lambda_i^{x_i} / x_i!)}{e^{-\lambda_+} \lambda_+^z / z!} \\ &= \binom{z}{x_1 \cdots x_n} \left(\frac{\lambda_1}{\lambda_+} \right)^{x_1} \cdots \left(\frac{\lambda_n}{\lambda_+} \right)^{x_n}. \end{aligned}$$

□

2.9.3 The distribution of the maximum of two random variables

Let X_1 and X_2 be independent and both have distribution F . Suppose that F has density f w.r.t. Lebesgue measure. Let

$$Z := \max\{X_1, X_2\}.$$

Lemma *The distribution function of Z is F^2 . Moreover, Z has density*

$$f_Z(z) = 2F(z)f(z), \quad z \in \mathbb{R}.$$

Proof. We have for all z ,

$$\begin{aligned} P(Z \leq z) &= P(\max\{X_1, X_2\} \leq z) \\ &= P(X_1 \leq z, X_2 \leq z) = F^2(z). \end{aligned}$$

If F has density f , then (Lebesgue)-almost everywhere,

$$f(z) = \frac{d}{dz}F(z).$$

So the derivative of F^2 exists almost everywhere, and

$$\frac{d}{dz}F^2(z) = 2F(z)f(z).$$

□

Let $X := (X_1, X_2)$. The conditional density of X given $Z = z$ has density

$$f_X(x_1, x_2|z) = \begin{cases} \frac{f(x_2)}{2F(z)} & \text{if } x_1 = z \text{ and } x_2 < z \\ \frac{f(x_1)}{2F(z)} & \text{if } x_1 < z \text{ and } x_2 = z \\ 0 & \text{else} \end{cases}.$$

The conditional distribution function of X_1 given $Z = z$ is

$$F_{X_1}(x_1|z) = \begin{cases} \frac{F(x_1)}{2F(z)}, & x_1 < z \\ 1, & x_1 \geq z \end{cases}.$$

Note thus that this distribution has a jump of size $1/2$ at z .

2.10 Sufficiency

Let $S : \mathcal{X} \rightarrow \mathcal{Y}$ be some given map. We consider the statistic $S = S(X)$. Throughout, by the phrase *for all possible s* , we mean for all s for which conditional distributions given $S = s$ are defined (in other words: for all s in the support of the distribution of S , which may depend on θ).

Definition *We call S sufficient for $\theta \in \Theta$ if for all θ , and all possible s , the conditional distribution*

$$P_\theta(X \in \cdot | S(X) = s)$$

does not depend on θ .

Example 2.10.1 Let X_1, \dots, X_n be i.i.d., and have the Bernoulli distribution with probability $\theta \in (0, 1)$ of success: (for $i = 1, \dots, n$)

$$P_\theta(X_i = 1) = 1 - P_\theta(X_i = 0) = \theta.$$

Take $S = \sum_{i=1}^n X_i$. Then S is sufficient for θ : for all possible s ,

$$\mathbb{P}_\theta(X_1 = x_1, \dots, X_n = x_n | S = s) = \frac{1}{\binom{n}{s}}, \sum_{i=1}^n x_i = s.$$

Example 2.10.2 Let $\mathbf{X} := (X_1, \dots, X_n)$, with X_1, \dots, X_n i.i.d. and Poisson(θ)-distributed. Take $S = \sum_{i=1}^n X_i$. Then S has the Poisson($n\theta$)-distribution. For all possible s , the conditional distribution of \mathbf{X} given $S = s$ is the multinomial distribution with parameters s and $(p_1, \dots, p_n) = (\frac{1}{n}, \dots, \frac{1}{n})$:

$$\mathbb{P}_\theta(X_1 = x_1, \dots, X_n = x_n | S = s) = \binom{s}{x_1 \dots x_n} \left(\frac{1}{n}\right)^s, \sum_{i=1}^n x_i = s.$$

This distribution does not depend on θ , so S is sufficient for θ .

Example 2.10.3 Let X_1 and X_2 be independent, and both have the exponential distribution with parameter $\theta > 0$. The density of e.g., X_1 is then

$$f_{X_1}(x; \theta) = \theta e^{-\theta x}, \quad x > 0.$$

Let $S = X_1 + X_2$. Verify that S has density

$$f_S(s; \theta) = s\theta^2 e^{-\theta s}, \quad s > 0.$$

(This is the Gamma(2, θ)-distribution.) For all possible s , the conditional density of (X_1, X_2) given $S = s$ is thus

$$f_{X_1, X_2}(x_1, x_2 | S = s) = \frac{1}{s}, \quad x_1 + x_2 = s.$$

Hence, S is sufficient for θ .

Example 2.10.4 Let X_1, \dots, X_n be an i.i.d. sample from a continuous distribution F . Then $S := (X_{(1)}, \dots, X_{(n)})$ is sufficient for F : for all possible $s = (s_1, \dots, s_n)$ ($s_1 < \dots < s_n$), and for $(x_{q_1}, \dots, x_{q_n}) = s$,

$$\mathbb{P}_\theta \left((X_1, \dots, X_n) = (x_1, \dots, x_n) \mid (X_{(1)}, \dots, X_{(n)}) = s \right) = \frac{1}{n!}.$$

Example 2.10.5 Let X_1 and X_2 be independent, and both uniformly distributed on the interval $[0, \theta]$, with $\theta > 0$. Define $Z := X_1 + X_2$.

Lemma *The random variable Z has density*

$$f_Z(z; \theta) = \begin{cases} z/\theta^2 & \text{if } 0 \leq z \leq \theta \\ (2\theta - z)/\theta^2 & \text{if } \theta \leq z \leq 2\theta \end{cases}.$$

Proof. First, assume $\theta = 1$. Then the distribution function of Z is

$$F_Z(z) = \begin{cases} z^2/2 & 0 \leq z \leq 1 \\ 1 - (2-z)^2/2 & 1 \leq z \leq 2 \end{cases} .$$

So the density is then

$$f_Z(z) = \begin{cases} z & 0 \leq z \leq 1 \\ 2-z & 1 \leq z \leq 2 \end{cases} .$$

For general θ , the result follows from the uniform case by the transformation $Z \mapsto \theta Z$, which maps f_Z into $f_Z(\cdot/\theta)/\theta$. \square

The conditional density of (X_1, X_2) given $Z = z \in (0, 2\theta)$ is now

$$f_{X_1, X_2}(x_1, x_2 | Z = z; \theta) = \begin{cases} \frac{1}{z} & 0 \leq z \leq \theta \\ \frac{1}{2\theta-z} & \theta \leq z \leq 2\theta \end{cases} .$$

This depends on θ , so Z is not sufficient for θ .

Consider now $S := \max\{X_1, X_2\}$. The conditional density of (X_1, X_2) given $S = s \in (0, \theta)$ is

$$f_{X_1, X_2}(x_1, x_2 | S = s) = \frac{1}{2s}, \quad 0 \leq x_1 < s, \quad x_2 = s \text{ or } x_1 = s, \quad 0 \leq x_2 < s.$$

This does not depend on θ , so S is sufficient for θ .

2.10.1 Rao-Blackwell

Lemma 2.10.1 *Suppose S is sufficient for θ . Let $d : \mathcal{X} \rightarrow \mathcal{A}$ be some decision. Then there is a randomized decision $\delta(S)$ that only depends on S , such that*

$$R(\theta, \delta(S)) = R(\theta, d), \quad \forall \theta.$$

Proof. Let X_s^* be a random variable with distribution $P(X \in \cdot | S = s)$. Then, by construction, for all possible s , the conditional distribution, given $S = s$, of X_s^* and X are equal. It follows that X and X_s^* have the same distribution. Formally, let us write Q_θ for the distribution of S . Then

$$\begin{aligned} P_\theta(X_s^* \in \cdot) &= \int P(X_s^* \in \cdot | S = s) dQ_\theta(s) \\ &= \int P(X \in \cdot | S = s) dQ_\theta(s) = P_\theta(X \in \cdot). \end{aligned}$$

The result of the lemma follows by taking $\delta(s) := d(X_s^*)$. \square

Lemma 2.10.2 (Rao Blackwell) *Suppose that S is sufficient for θ . Suppose moreover that the action space $\mathcal{A} \subset \mathbb{R}^p$ is convex, and that for each θ , the map $a \mapsto L(\theta, a)$ is convex. Let $d : \mathcal{X} \rightarrow \mathcal{A}$ be a decision, and define $d'(s) := E(d(X) | S = s)$ (assumed to exist). Then*

$$R(\theta, d') \leq R(\theta, d), \quad \forall \theta.$$

Proof. Jensen's inequality says that for a convex function g ,

$$E(g(X)) \geq g(EX).$$

Hence, $\forall \theta$,

$$\begin{aligned} E\left(L\left(\theta, d(X)\right) \middle| S = s\right) &\geq L\left(\theta, E\left(d(X) \middle| S = s\right)\right) \\ &= L(\theta, d'(s)). \end{aligned}$$

By the iterated expectations lemma, we arrive at

$$\begin{aligned} R(\theta, d) &= E_{\theta}L(\theta, d(X)) \\ &= E_{\theta}E\left(L\left(\theta, d(X)\right) \middle| S\right) \geq E_{\theta}L(\theta, d'(S)). \end{aligned}$$

□

2.10.2 Factorization Theorem of Neyman

Theorem 2.10.1 (*Factorization Theorem of Neyman*) Suppose $\{P_{\theta} : \theta \in \Theta\}$ is dominated by a σ -finite measure ν . Let $p_{\theta} := dP_{\theta}/d\nu$ denote the densities. Then S is sufficient for θ if and only if one can write p_{θ} in the form

$$p_{\theta}(x) = g_{\theta}(S(x))h(x), \quad \forall x, \theta.$$

Proof in the discrete case. Suppose X takes only the values $a_1, a_2, \dots \forall \theta$ (so we may take ν to be the counting measure). Let Q_{θ} be the distribution of S :

$$Q_{\theta}(s) := \sum_{j: S(a_j)=s} P_{\theta}(X = a_j).$$

The conditional distribution of X given S is

$$P_{\theta}(X = x | S = s) = \frac{P_{\theta}(X = x)}{Q_{\theta}(s)}, \quad S(x) = s.$$

(\Rightarrow) If S is sufficient for θ , the above does not depend on θ , but is only a function of x , say $h(x)$. So we may write for $S(x) = s$,

$$P_{\theta}(X = x) = P_{\theta}(X = x | S = s)Q_{\theta}(S = s) = h(x)g_{\theta}(s),$$

with $g_{\theta}(s) = Q_{\theta}(S = s)$.

(\Leftarrow) Inserting $p_{\theta}(x) = g_{\theta}(S(x))h(x)$, we find

$$Q_{\theta}(s) = g_{\theta}(s) \sum_{j: S(a_j)=s} h(a_j),$$

This gives in the formula for $P_\theta(X = x|S = s)$,

$$P_\theta(X = x|S = s) = \frac{h(x)}{\sum_{j: S(a_j)=s} h(a_j)}$$

which does not depend on θ . □

Remark The proof for the general case is along the same lines, but does have some subtle elements!



Corollary 2.10.1 *The likelihood is $L_X(\theta) = p_\theta(X) = g_\theta(S)h(X)$. Hence, the maximum likelihood estimator $\hat{\theta} = \arg \max_\theta L_X(\theta) = \arg \max_\theta g_\theta(S)$ depends only on the sufficient statistic S .*

Corollary 2.10.2 *The Bayes decision is*

$$d_{\text{Bayes}}(X) = \arg \min_{a \in \mathcal{A}} l(X, a),$$

where

$$\begin{aligned} l(x, a) &= E(L(\theta, a)|X = x) = \int L(\vartheta, a)w(\vartheta|x)d\mu(\vartheta) \\ &= \int L(\vartheta, a)g_\vartheta(S(x))w(\vartheta)d\mu(\vartheta)h(x)/p(x). \end{aligned}$$

So

$$d_{\text{Bayes}}(X) = \arg \min_{a \in \mathcal{A}} \int L(\vartheta, a)g_\vartheta(S)w(\vartheta)d\mu(\vartheta),$$

which only depends on the sufficient statistic S .

Example 2.10.6 Let X_1, \dots, X_n be i.i.d., and uniformly distributed on the interval $[0, \theta]$. Then the density of $\mathbf{X} = (X_1, \dots, X_n)$ is

$$\begin{aligned} p_\theta(x_1, \dots, x_n) &= \frac{1}{\theta^n} \mathbf{1}\{0 \leq \min\{x_1, \dots, x_n\} \leq \max\{x_1, \dots, x_n\} \leq \theta\} \\ &= g_\theta(S(x_1, \dots, x_n))h(x_1, \dots, x_n), \end{aligned}$$

with

$$g_\theta(s) := \frac{1}{\theta^n} \mathbf{1}\{s \leq \theta\},$$

and

$$h(x_1, \dots, x_n) := \mathbf{1}\{0 \leq \min\{x_1, \dots, x_n\}\}.$$

Thus, $S = \max\{X_1, \dots, X_n\}$ is sufficient for θ .

2.10.3 Exponential families

Definition A k -dimensional exponential family is a family of distributions $\{P_\theta : \theta \in \Theta\}$, dominated by some σ -finite measure ν , with densities $p_\theta = dP_\theta/d\nu$ of the form

$$p_\theta(x) = \exp\left[\sum_{j=1}^k c_j(\theta)T_j(x) - d(\theta)\right]h(x).$$

Note In case of a k -dimensional exponential family, the k -dimensional statistic $S(X) = (T_1(X), \dots, T_k(X))$ is sufficient for θ .

Note If X_1, \dots, X_n is an i.i.d. sample from a k -dimensional exponential family, then the distribution of $\mathbf{X} = (X_1, \dots, X_n)$ is also in a k -dimensional exponential family. The density of \mathbf{X} is then (for $\mathbf{x} := (x_1, \dots, x_n)$),

$$\mathbf{p}_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i) = \exp\left[\sum_{j=1}^k nc_j(\theta)\bar{T}_j(\mathbf{x}) - nd(\theta)\right] \prod_{i=1}^n h(x_i),$$

where, for $j = 1, \dots, k$,

$$\bar{T}_j(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n T_j(x_i).$$

Hence $S(\mathbf{X}) = (\bar{T}_1(\mathbf{X}), \dots, \bar{T}_k(\mathbf{X}))$ is then sufficient for θ .

Note The functions $\{T_j\}$ and $\{c_j\}$ are not uniquely defined.

Example 2.10.7 If X is Poisson(θ)-distributed, we have

$$\begin{aligned} p_\theta(x) &= e^{-\theta} \frac{\theta^x}{x!} \\ &= \exp[x \log \theta - \theta] \frac{1}{x!}. \end{aligned}$$

Hence, we may take $T(x) = x$, $c(\theta) = \log \theta$, and $d(\theta) = \theta$.

Example 2.10.8 If X has the Binomial(n, θ)-distribution, we have

$$\begin{aligned} p_\theta(x) &= \binom{n}{x} \theta^x (1 - \theta)^{n-x} \\ &= \binom{n}{x} \left(\frac{\theta}{1 - \theta}\right)^x (1 - \theta)^n \\ &= \binom{n}{x} \exp\left[x \log\left(\frac{\theta}{1 - \theta}\right) + n \log(1 - \theta)\right]. \end{aligned}$$

So we can take $T(x) = x$, $c(\theta) = \log(\theta/(1 - \theta))$, and $d(\theta) = -n \log(1 - \theta)$.

Example 2.10.9 If X has the Negative Binomial(k, θ)-distribution we have

$$\begin{aligned} p_\theta(x) &= \frac{\Gamma(x+k)}{\Gamma(k)x!} \theta^k (1-\theta)^x \\ &= \frac{\Gamma(x+k)}{\Gamma(k)x!} \exp[x \log(1-\theta) + k \log(\theta)]. \end{aligned}$$

So we may take $T(x) = x$, $c(\theta) = \log(1-\theta)$ and $d(\theta) = -k \log(\theta)$.

Example 2.10.10 Let X have the Gamma(k, θ)-distribution (with k known). Then

$$\begin{aligned} p_\theta(x) &= e^{-\theta x} x^{k-1} \frac{\theta^k}{\Gamma(k)} \\ &= \frac{x^{k-1}}{\Gamma(k)} \exp[-\theta x + k \log \theta]. \end{aligned}$$

So we can take $T(x) = x$, $c(\theta) = -\theta$, and $d(\theta) = -k \log \theta$.

Example 2.10.11 Let X have the Gamma(k, λ)-distribution, and let $\theta = (k, \lambda)$. Then

$$\begin{aligned} p_\theta(x) &= e^{-\lambda x} x^{k-1} \frac{\lambda^k}{\Gamma(k)} \\ &= \frac{x^{k-1}}{\Gamma(k)} \exp[-\lambda x + (k-1) \log x + k \log \lambda - \log \Gamma(k)]. \end{aligned}$$

So we can take $T_1(x) = x$, $T_2(x) = \log x$, $c(\theta) = -\lambda$, $c_2(\theta) = (k-1)$, and $d(\theta) = -k \log \lambda + \log \Gamma(k)$.

Example 2.10.12 Let X be $\mathcal{N}(\mu, \sigma^2)$ -distributed, and let $\theta = (\mu, \sigma)$. Then

$$\begin{aligned} p_\theta(x) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right] \\ &= \frac{1}{\sqrt{2\pi}} \exp\left[\frac{x\mu}{\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} - \log \sigma\right]. \end{aligned}$$

So we can take $T_1(x) = x$, $T_2(x) = x^2$, $c_1(\theta) = \mu/\sigma^2$, $c_2(\theta) = -1/(2\sigma^2)$, and $d(\theta) = \mu^2/(2\sigma^2) + \log(\sigma)$.

2.10.4 Canonical form of an exponential family

In this subsection, we assume regularity conditions, such as existence of derivatives, and inverses, and permission to interchange differentiation and integration.



Let $\Theta \subset \mathbb{R}^k$, and let $\{P_\theta : \theta \in \Theta\}$ be a family of probability measures dominated by a σ -finite measure ν . Define the densities

$$p_\theta := \frac{dP_\theta}{d\nu}.$$

Definition We call $\{P_\theta : \theta \in \Theta\}$ an exponential family in canonical form, if

$$p_\theta(x) = \exp \left[\sum_{j=1}^k \theta_j T_j(x) - d(\theta) \right] h(x).$$

Note that $d(\theta)$ is the normalizing constant

$$d(\theta) = \log \left(\int \exp \left[\sum_{j=1}^k \theta_j T_j(x) \right] h(x) d\nu(x) \right).$$

We let

$$\dot{d}(\theta) := \frac{\partial}{\partial \theta} d(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} d(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_k} d(\theta) \end{pmatrix}$$

denote the vector of first derivatives, and

$$\ddot{d}(\theta) := \frac{\partial^2}{\partial \theta \partial \theta^\top} d(\theta) = \left(\frac{\partial^2 d(\theta)}{\partial \theta_j \partial \theta_{j'}} \right)$$

denote the $k \times k$ matrix of second derivatives. Further, we write

$$T(X) := \begin{pmatrix} T_1(X) \\ \vdots \\ T_k(X) \end{pmatrix}, \quad E_\theta T(X) := \begin{pmatrix} E_\theta T_1(X) \\ \vdots \\ E_\theta T_k(X) \end{pmatrix},$$

and we write the $k \times k$ covariance matrix of $T(X)$ as

$$\text{Cov}_\theta(T(X)) := \left(\text{cov}_\theta(T_j(X), T_{j'}(X)) \right).$$

Lemma We have (under regularity)

$$E_\theta T(X) = \dot{d}(\theta), \quad \text{Cov}_\theta(T(X)) = \ddot{d}(\theta).$$

Proof. By the definition of $d(\theta)$, we find

$$\begin{aligned} \dot{d}(\theta) &= \frac{\partial}{\partial \theta} \log \left(\int \exp \left[\theta^\top T(x) \right] h(x) d\nu(x) \right) \\ &= \frac{\int \exp \left[\theta^\top T(x) \right] T(x) h(x) d\nu(x)}{\int \exp \left[\theta^\top T(x) \right] h(x) d\nu(x)} \\ &= \int \exp \left[\theta^\top T(x) - d(\theta) \right] T(x) h(x) d\nu(x) \end{aligned}$$

$$= \int p_\theta(x)T(x)d\nu(x) = E_\theta T(X),$$

and

$$\begin{aligned} \ddot{d}(\theta) &= \frac{\int \exp\left[\theta^\top T(x)\right]T(x)T(x)^\top h(x)d\nu(x)}{\int \exp\left[\theta^\top T(x)\right]h(x)d\nu(x)} \\ &= \frac{\left(\int \exp\left[\theta^\top T(x)\right]T(x)h(x)d\nu(x)\right)\left(\int \exp\left[\theta^\top T(x)\right]T(x)h(x)d\nu(x)\right)^\top}{\left(\int \exp\left[\theta^\top T(x)\right]h(x)d\nu(x)\right)^2} \\ &= \int \exp\left[\theta^\top T(x) - d(\theta)\right]T(x)T(x)^\top h(x)d\nu(x) \\ &\quad - \left(\int \exp\left[\theta^\top T(x) - d(\theta)\right]T(x)h(x)d\nu(x)\right) \\ &\quad \times \left(\int \exp\left[\theta^\top T(x) - d(\theta)\right]T(x)h(x)d\nu(x)\right)^\top \\ &= \int p_\theta(x)T(x)T(x)^\top d\nu(x) \\ &\quad - \left(\int p_\theta(x)T(x)d\nu(x)\right)\left(\int p_\theta(x)T(x)d\nu(x)\right)^\top \\ &= E_\theta T(X)T(X)^\top - \left(E_\theta T(X)\right)\left(E_\theta T(X)\right)^\top \\ &= \text{Cov}_\theta(T(X)). \end{aligned}$$

□

Let us now simplify to the one-dimensional case, that is $\Theta \subset \mathbb{R}$. Consider an exponential family, not necessarily in canonical form:

$$p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x).$$

We can put this in canonical form by reparametrizing

$$\theta \mapsto c(\theta) := \gamma \text{ (say),}$$

to get

$$\tilde{p}_\gamma(x) = \exp[\gamma T(x) - d_0(\gamma)]h(x),$$

where

$$d_0(\gamma) = d(c^{-1}(\gamma)).$$

It follows that

$$E_\theta T(X) = \dot{d}_0(\gamma) = \frac{\dot{d}(c^{-1}(\gamma))}{\dot{c}(c^{-1}(\gamma))} = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}, \quad (2.2)$$

and

$$\begin{aligned} \text{var}_\theta(T(X)) &= \ddot{d}_0(\gamma) = \frac{\ddot{d}(c^{-1}(\gamma))}{[\dot{c}(c^{-1}(\gamma))]^2} - \frac{\dot{d}(c^{-1}(\gamma))\ddot{c}(c^{-1}(\gamma))}{[\dot{c}(c^{-1}(\gamma))]^3} \\ &= \frac{\ddot{d}(\theta)}{[\dot{c}(\theta)]^2} - \frac{\dot{d}(\theta)\ddot{c}(\theta)}{[\dot{c}(\theta)]^3} = \frac{1}{[\dot{c}(\theta)]^2} \left(\ddot{d}(\theta) - \frac{\dot{d}(\theta)}{\dot{c}(\theta)}\ddot{c}(\theta) \right). \end{aligned} \quad (2.3)$$

For an arbitrary (but regular) family of densities $\{p_\theta : \theta \in \Theta\}$, with (again for simplicity) $\Theta \subset \mathbb{R}$, we define the *score function*

$$s_\theta(x) := \frac{d}{d\theta} \log p_\theta(x),$$

and the *Fisher information* for estimating θ

$$I(\theta) := \text{var}_\theta(s_\theta(X))$$

(see also Chapter 3 and 6).

Lemma *We have (under regularity)*

$$E_\theta s_\theta(X) = 0,$$

and

$$I(\theta) = -E_\theta \dot{s}_\theta(X),$$

where $\dot{s}_\theta(x) := \frac{d}{d\theta} s_\theta(x)$.

Proof. The results follow from the fact that densities integrate to one, assuming that we may interchange derivatives and integrals:

$$\begin{aligned} E_\theta s_\theta(X) &= \int s_\theta(x) p_\theta(x) d\nu(x) \\ &= \int \frac{d \log p_\theta(x)}{d\theta} p_\theta(x) d\nu(x) = \int \frac{dp_\theta(x)/d\theta}{p_\theta(x)} p_\theta(x) d\nu(x) \\ &= \int \frac{d}{d\theta} p_\theta(x) d\nu(x) = \frac{d}{d\theta} \int p_\theta(x) d\nu(x) = \frac{d}{d\theta} 1 = 0, \end{aligned}$$

and

$$\begin{aligned} E_\theta \dot{s}_\theta(X) &= E_\theta \left[\frac{d^2 p_\theta(X)/d\theta^2}{p_\theta(X)} - \left(\frac{dp_\theta(X)/d\theta}{p_\theta(X)} \right)^2 \right] \\ &= E_\theta \left[\frac{d^2 p_\theta(X)/d\theta^2}{p_\theta(X)} \right] - E_\theta s_\theta^2(X). \end{aligned}$$

Now, $E_\theta s_\theta^2(X)$ equals $\text{var}_\theta s_\theta(X)$, since $E_\theta s_\theta(X) = 0$. Moreover,

$$E_\theta \left[\frac{d^2 p_\theta(X)/d\theta^2}{p_\theta(X)} \right] = \int \frac{d^2}{d\theta^2} p_\theta(x) d\nu(x) = \frac{d^2}{d\theta^2} \int p_\theta(x) d\nu(x) = \frac{d^2}{d\theta^2} 1 = 0.$$

□

In the special case that $\{P_\theta : \theta \in \Theta\}$ is a one-dimensional exponential family, the densities are of the form

$$p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x).$$

Hence

$$s_\theta(x) = \dot{c}(\theta)T(x) - \dot{d}(\theta).$$

The equality $E_\theta s_\theta(X) = 0$ implies that

$$E_\theta T(X) = \frac{\dot{d}(\theta)}{\dot{c}(\theta)},$$

which re-establishes (2.2). One moreover has

$$\dot{s}_\theta(x) = \ddot{c}(\theta)T(x) - \ddot{d}(\theta).$$

Hence, the inequality $\text{var}_\theta(s_\theta(X)) = -E_\theta \dot{s}_\theta(X)$ implies

$$\begin{aligned} [\dot{c}(\theta)]^2 \text{var}_\theta(T(X)) &= -\ddot{c}(\theta)E_\theta T(X) + \ddot{d}(\theta) \\ &= \ddot{d}(\theta) - \frac{\dot{d}(\theta)}{\dot{c}(\theta)}\ddot{c}(\theta), \end{aligned}$$

which re-establishes (2.3). In addition, it follows that

$$I(\theta) = \ddot{d}(\theta) - \frac{\dot{d}(\theta)}{\dot{c}(\theta)}\ddot{c}(\theta).$$

The Fisher information for estimating $\gamma = c(\theta)$ is

$$I_0(\gamma) = \ddot{d}_0(\gamma) = \frac{I(\theta)}{[\dot{c}(\theta)]^2}.$$

More generally, the Fisher information for estimating a differentiable function $g(\theta)$ of the parameter θ , is equal to $I(\theta)/[\dot{g}(\theta)]^2$.

Example

Let $X \in \{0, 1\}$ have the Bernoulli-distribution with success parameter $\theta \in (0, 1)$:

$$p_\theta(x) = \theta^x(1 - \theta)^{1-x} = \exp\left[x \log\left(\frac{\theta}{1 - \theta}\right) + \log(1 - \theta)\right], \quad x \in \{0, 1\}.$$

We reparametrize:

$$\gamma := c(\theta) = \log\left(\frac{\theta}{1 - \theta}\right),$$

which is called the log-odds ratio. Inverting gives

$$\theta = \frac{e^\gamma}{1 + e^\gamma},$$

and hence

$$d(\theta) = -\log(1 - \theta) = \log\left(1 + e^\gamma\right) := d_0(\gamma).$$

Thus

$$\dot{d}_0(\gamma) = \frac{e^\gamma}{1 + e^\gamma} = \theta = E_\theta X,$$

and

$$\ddot{d}_0(\gamma) = \frac{e^\gamma}{1 + e^\gamma} - \frac{e^{2\gamma}}{(1 + e^\gamma)^2} = \frac{e^\gamma}{(1 + e^\gamma)^2} = \theta(1 - \theta) = \text{var}_\theta(X).$$

The score function is

$$\begin{aligned} s_\theta(x) &= \frac{d}{d\theta} \left[x \log \left(\frac{\theta}{1 - \theta} \right) + \log(1 - \theta) \right] \\ &= \frac{x}{\theta(1 - \theta)} - \frac{1}{1 - \theta}. \end{aligned}$$

The Fisher information for estimating the success parameter θ is

$$E_\theta s_\theta^2(X) = \frac{\text{var}_\theta(X)}{[\theta(1 - \theta)]^2} = \frac{1}{\theta(1 - \theta)},$$

whereas the Fisher information for estimating the log-odds ratio γ is

$$I_0(\gamma) = \theta(1 - \theta).$$

2.10.5 Minimal sufficiency

Definition We say that two likelihoods $L_x(\theta)$ and $L_{x'}(\theta)$ are proportional at (x, x') , if

$$L_x(\theta) = L_{x'}(\theta)c(x, x'), \forall \theta,$$

for some constant $c(x, x')$.

A statistic S is called minimal sufficient if $S(x) = S(x')$ for all x and x' for which the likelihoods are proportional.



Example 2.10.13 Let X_1, \dots, X_n be independent and $\mathcal{N}(\theta, 1)$ -distributed. Then $S = \sum_{i=1}^n X_i$ is sufficient for θ . We moreover have

$$\log L_{\mathbf{x}}(\theta) = S(\mathbf{x})\theta - \frac{n\theta^2}{2} - \frac{\sum_{i=1}^n x_i^2}{2} - \log(2\pi)/2.$$

So

$$\log L_{\mathbf{x}}(\theta) - \log L_{\mathbf{x}'}(\theta) = (S(\mathbf{x}) - S(\mathbf{x}'))\theta - \frac{\sum_{i=1}^n x_i^2 - \sum_{i=1}^n (x'_i)^2}{2},$$

which equals,

$$\log c(\mathbf{x}, \mathbf{x}'), \forall \theta,$$

for some function c , if and only if $S(\mathbf{x}) = S(\mathbf{x}')$. So S is minimal sufficient.

Example 2.10.14 Let X_1, \dots, X_n be independent and Laplace-distributed with location parameter θ . Then

$$\log L_{\mathbf{x}}(\theta) = -(\log 2)/2 - \sqrt{2} \sum_{i=1}^n |x_i - \theta|,$$

so

$$\log L_{\mathbf{x}}(\theta) - \log L_{\mathbf{x}'}(\theta) = -\sqrt{2} \sum_{i=1}^n (|x_i - \theta| - |x'_i - \theta|)$$

which equals

$$\log c(\mathbf{x}, \mathbf{x}'), \quad \forall \theta,$$

for some function c , if and only if $(x_{(1)}, \dots, x_{(n)}) = (x'_{(1)}, \dots, x'_{(n)})$. So the order statistics $X_{(1)}, \dots, X_{(n)}$ are minimal sufficient.

Chapter 3

Unbiased estimators

3.1 What is an unbiased estimator?

Let $X \in \mathcal{X}$ denote the observations. The distribution P of X is assumed to be a member of a given class $\{P_\theta : \theta \in \Theta\}$ of distributions. The parameter of interest in this chapter is $\gamma := g(\theta)$, with $g : \Theta \rightarrow \mathbb{R}$ (for simplicity, we initially assume γ to be one-dimensional).

Let $T : \mathcal{X} \rightarrow \mathbb{R}$ be an estimator of $g(\theta)$.

Definition *The bias of $T = T(X)$ is*

$$\text{bias}_\theta(T) := E_\theta T - g(\theta).$$

The estimator T is called unbiased if

$$\text{bias}_\theta(T) = 0, \forall \theta.$$

Thus, unbiasedness means that there is no systematic error: $E_\theta T = g(\theta)$. We require this **for all** θ .

Example 3.1.1 Let $X \sim \text{Binomial}(n, \theta)$, $0 < \theta < 1$. We have

$$E_\theta T(X) = \sum_{k=0}^n \binom{n}{k} \theta^k (1-\theta)^{n-k} T(k) := q(\theta).$$

Note that $q(\theta)$ is a polynomial in θ of degree at most n . So only parameters $g(\theta)$ which are polynomials of degree at most n can be estimated unbiasedly. It means that there exists no unbiased estimator of, for example, $\sqrt{\theta}$ or $\theta/(1-\theta)$.

Example 3.1.2 Let $X \sim \text{Poisson}(\theta)$. Then

$$E_\theta T(X) = \sum_{k=0}^{\infty} e^{-\theta} \frac{\theta^k}{k!} T(k) := e^{-\theta} p(\theta).$$

Note that $p(\theta)$ is a power series in θ . Thus only parameters $g(\theta)$ which are a power series in θ times $e^{-\theta}$ can be estimated unbiasedly. An example is the probability of early failure

$$g(\theta) := e^{-\theta} = P_{\theta}(X = 0).$$

An unbiased estimator of $e^{-\theta}$ is for instance

$$T(X) = 1\{X = 0\}.$$

As another example, suppose the parameter of interest is

$$g(\theta) := e^{-2\theta}.$$

An unbiased estimator is

$$T(X) = \begin{cases} +1 & \text{if } X \text{ is even} \\ -1 & \text{if } X \text{ is odd} \end{cases}.$$

This estimator does not make sense at all!

Example 3.1.3 Let X_1, \dots, X_n be i.i.d. $\mathcal{N}(\mu, \sigma^2)$, and let $\theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+$. Then

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

is an unbiased estimator of σ^2 . But S is not an unbiased estimator of σ . In fact, one can show that there does not exist any unbiased estimator of σ !

We conclude that requiring unbiasedness can have disadvantages: unbiased estimators do not always exist, and if they do, they can be nonsensical. Moreover, the property of unbiasedness is not preserved under taking nonlinear transformations.

3.2 UMVU estimators

Lemma 3.2.1 *We have the following equality for the mean square error:*

$$E_{\theta}|T - g(\theta)|^2 = \text{bias}_{\theta}^2(T) + \text{var}_{\theta}(T).$$

In other words, the mean square error consists of two components, the (squared) bias and the variance. This is called the bias-variance decomposition. As we will see, it is often the case that an attempt to decrease the bias results in an increase of the variance (and vice versa).

Example 3.2.1 Let X_1, \dots, X_n be i.i.d. $\mathcal{N}(\mu, \sigma^2)$ -distributed. Both μ and σ^2 are unknown parameters: $\theta := (\mu, \sigma^2)$.

Case i Suppose the mean μ is our parameter of interest. Consider the estimator $T := \lambda\bar{X}$, where $0 \leq \lambda \leq 1$. Then the bias is decreasing in λ , but the variance is increasing in λ :

$$E_{\theta}|T - \mu|^2 = (1 - \lambda)^2\mu^2 + \lambda^2\sigma^2/n.$$

The right hand side can be minimized as a function of λ . The minimum is attained at

$$\lambda_{\text{opt}} := \frac{\mu^2}{\sigma^2/n + \mu^2}.$$

However, $\lambda_{\text{opt}}\bar{X}$ is not an estimator as it depends on the unknown parameters.

Case ii Suppose σ^2 is the parameter of interest. Let S^2 be the sample variance:

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

It is known that S^2 is unbiased. But does it also have small mean square error? Let us compare it with the estimator

$$\hat{\sigma}^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

To compute the mean square errors of these two estimators, we first recall that

$$\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2,$$

a χ^2 -distribution with $n-1$ degrees of freedom. The χ^2 -distribution is a special case of the Gamma-distribution, namely

$$\chi_{n-1}^2 = \Gamma\left(\frac{n-1}{2}, \frac{1}{2}\right).$$

Thus ¹

$$E_{\theta} \left(\sum_{i=1}^n (X_i - \bar{X})^2 / \sigma^2 \right) = n-1, \quad \text{var} \left(\sum_{i=1}^n (X_i - \bar{X})^2 / \sigma^2 \right) = 2(n-1).$$

It follows that

$$E_{\theta}|S^2 - \sigma^2|^2 = \text{var}(S^2) = \frac{\sigma^4}{(n-1)^2} 2(n-1) = \frac{2\sigma^4}{n-1},$$

and

$$E_{\theta}\hat{\sigma}^2 = \frac{n-1}{n}\sigma^2,$$

$$\text{bias}_{\theta}(\hat{\sigma}^2) = -\frac{1}{n}\sigma^2,$$

¹If Y has a $\Gamma(k, \lambda)$ -distribution, then $EY = k/\lambda$ and $\text{var}(Y) = k/\lambda^2$.

so that

$$E_{\theta}|\hat{\sigma}^2 - \sigma^2|^2 = \text{bias}_{\theta}^2(\hat{\sigma}^2) + \text{var}_{\theta}(\hat{\sigma}^2) = \frac{\sigma^4}{n^2} + \frac{\sigma^4}{n^2}2(n-1) = \frac{\sigma^4(2n-1)}{n^2}.$$

Conclusion: the mean square error of $\hat{\sigma}^2$ is smaller than the mean square error of S^2 !

Generally, it is not possible to construct an estimator that possesses the best among all of all desirable properties. We therefore fix one property: unbiasedness (despite its disadvantages), and look for good estimators among the unbiased ones.

Definition An unbiased estimator T^* is called UMVU (Uniform Minimum Variance Unbiased) if for any other unbiased estimator T ,

$$\text{var}_{\theta}(T^*) \leq \text{var}_{\theta}(T), \quad \forall \theta.$$

Suppose that T is unbiased, and that S is sufficient. Let

$$T^* := E(T|S).$$

The distribution of T given S does not depend on θ , so T^* is also an estimator. Moreover, it is unbiased:

$$E_{\theta}T^* = E_{\theta}(E(T|S)) = E_{\theta}T = g(\theta).$$

By conditioning on S , “superfluous” variance in the sample is killed. Indeed, the following lemma (which is a general property of conditional distributions) shows that T^* cannot have larger variance than T :

$$\text{var}_{\theta}(T^*) \leq \text{var}_{\theta}(T), \quad \forall \theta.$$

Lemma 3.2.2 Let Y and Z be two random variables. Then

$$\text{var}(Y) = \text{var}(E(Y|Z)) + E\text{var}(Y|Z).$$

Proof. It holds that

$$\begin{aligned} \text{var}(E(Y|Z)) &= E\left[E(Y|Z)\right]^2 - \left[E(E(Y|Z))\right]^2 \\ &= E[E(Y|Z)]^2 - [EY]^2, \end{aligned}$$

and

$$\begin{aligned} E\text{var}(Y|Z) &= E\left[E(Y^2|Z) - [E(Y|Z)]^2\right] \\ &= EY^2 - E[E(Y|Z)]^2. \end{aligned}$$

Hence, when adding up, the term $E[E(Y|Z)]^2$ cancels out, and what is left over is exactly the variance

$$\text{var}(Y) = EY^2 - [EY]^2.$$

□

3.2.1 Complete statistics

The question arises: can we construct an unbiased estimator with even smaller variance than $T^* = E(T|S)$? Note that T^* depends on X only via $S = S(X)$, i.e., it depends only on the sufficient statistic. In our search for UMVU estimators, we may restrict our attention to estimators depending only on S . Thus, if there is only one unbiased estimator depending only on S , it has to be UMVU.

Definition A statistic S is called complete if we have the following implication:

$$E_\theta h(S) = 0 \quad \forall \theta \Rightarrow h(S) = 0, \quad P_\theta - a.s., \forall \theta.$$

Lemma 3.2.3 (Lehmann-Scheffé) Let T be an unbiased estimator of $g(\theta)$, with, for all θ , finite variance. Moreover, let S be sufficient and complete. Then $T^* := E(T|S)$ is UMVU.

Proof. We already noted that $T^* = T^*(S)$ is unbiased and that $\text{var}_\theta(T^*) \leq \text{var}_\theta(T) \quad \forall \theta$. If $T'(S)$ is another unbiased estimator of $g(\theta)$, we have

$$E_\theta(T(S) - T'(S)) = 0, \quad \forall \theta.$$

Because S is complete, this implies

$$T^* = T', \quad P_\theta - a.s.$$

□

To check whether a statistic is complete, one often needs somewhat sophisticated tools from analysis/integration theory. In the next two examples, we only sketch the proofs of completeness.

Example 3.2.2 Let X_1, \dots, X_n be i.i.d. Poisson(θ)-distributed. We want to estimate $g(\theta) := e^{-\theta}$, the probability of early failure. An unbiased estimator is

$$T(X_1, \dots, X_n) := 1\{X_1 = 0\}.$$

A sufficient statistic is

$$S := \sum_{i=1}^n X_i.$$

We now check whether S is complete. Its distribution is the Poisson($n\theta$)-distribution. We therefore have for any function h ,

$$E_\theta h(S) = \sum_{k=0}^{\infty} e^{-n\theta} \frac{(n\theta)^k}{k!} h(k).$$

The equation

$$E_\theta h(S) = 0 \quad \forall \theta,$$

thus implies

$$\sum_{k=0}^{\infty} \frac{(n\theta)^k}{k!} h(k) = 0 \quad \forall \theta.$$

Let f be a function with Taylor expansion at zero. Then

$$f(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!} f^{(k)}(0).$$

The left hand side can only be zero for all x if $f \equiv 0$, in which case also $f^{(k)}(0) = 0$ for all k . Thus $(h(k))$ takes the role of $f^{(k)}(0)$ and $n\theta$ the role of x , we conclude that $h(k) = 0$ for all k , i.e., that S is complete.

So we know from the Lehmann-Scheffé Lemma that $T^* := E(T|S)$ is UMVU. Now,

$$\begin{aligned} P(T = 1|S = s) &= P(X_1 = 0|S = s) \\ &= \frac{e^{-\theta} e^{-(n-1)\theta} [(n-1)\theta]^s / s!}{e^{-n\theta} (n\theta)^s / s!} = \left(\frac{n-1}{n} \right)^s. \end{aligned}$$

Hence

$$T^* = \left(\frac{n-1}{n} \right)^S$$

is UMVU.

Example 3.2.3 Let X_1, \dots, X_n be i.i.d. Uniform $[0, \theta]$ -distributed, and $g(\theta) := \theta$. We know that $S := \max\{X_1, \dots, X_n\}$ is sufficient. The distribution function of S is

$$F_S(s) = P_{\theta}(\max\{X_1, \dots, X_n\} \leq s) = \left(\frac{s}{\theta} \right)^n, \quad 0 \leq s \leq \theta.$$

Its density is thus

$$f_S(s) = \frac{ns^{n-1}}{\theta^n}, \quad 0 \leq s \leq \theta.$$

Hence, for any (measurable) function h ,

$$E_{\theta}h(S) = \int_0^{\theta} h(s) \frac{ns^{n-1}}{\theta^n} ds.$$

If

$$E_{\theta}h(S) = 0 \quad \forall \theta,$$

it must hold that

$$\int_0^{\theta} h(s) s^{n-1} ds = 0 \quad \forall \theta.$$



Differentiating w.r.t. θ gives

$$h(\theta)\theta^{n-1} = 0 \quad \forall \theta,$$

which implies $h \equiv 0$. So S is complete.

It remains to find a statistic T^* that depends only on S and that is unbiased. We have

$$E_\theta S = \int_0^\theta s \frac{ns^{n-1}}{\theta^n} ds = \frac{n}{n+1} \theta.$$

So S itself is not unbiased, it is too small. But this can be easily repaired: take

$$T^* = \frac{n+1}{n} S.$$

Then, by the Lehmann-Scheffé Lemma, T^* is UMVU.

In the case of an exponential family, completeness holds for a sufficient statistic if the parameter space is “of the same dimension” as the sufficient statistic. This is stated more formally in the following lemma. We omit the proof.

Lemma 3.2.4 *Let for $\theta \in \Theta$,*

$$p_\theta(x) = \exp \left[\sum_{j=1}^k c_j(\theta) T_j(x) - d(\theta) \right] h(x).$$

Consider the set

$$\mathcal{C} := \{(c_1(\theta), \dots, c_k(\theta)) : \theta \in \Theta\} \subset \mathbb{R}^k.$$

Suppose that \mathcal{C} is truly k -dimensional (that is, not of dimension smaller than k), i.e., it contains an open ball in \mathbb{R}^k . (Or an open cube $\prod_{j=1}^k (a_j, b_j)$.) Then $S := (T_1, \dots, T_k)$ is complete.

Example 3.2.4 Let X_1, \dots, X_n be i.i.d. with $\Gamma(k, \lambda)$ -distribution. Both k and λ are assumed to be unknown, so that $\theta := (k, \lambda)$. We moreover let $\Theta := \mathbb{R}_+^2$. The density f of the $\Gamma(k, \lambda)$ -distribution is

$$f(z) = \frac{\lambda^k}{\Gamma(k)} e^{-\lambda z} z^{k-1}, \quad z > 0.$$

Hence,

$$p_\theta(x) = \exp \left[-\lambda \sum_{i=1}^n x_i + (k-1) \sum_{i=1}^n \log x_i - d(\theta) \right] h(x),$$

where

$$d(k, \lambda) = -nk \log \lambda + n \log \Gamma(k),$$

and

$$h(x) = \mathbb{1}\{x_i > 0, i = 1, \dots, n\}.$$

It follows that

$$\left(\sum_{i=1}^n X_i, \sum_{i=1}^n \log X_i \right)$$

is sufficient and complete.

Example 3.2.5 Consider two independent samples from normal distributions: X_1, \dots, X_n i.i.d. $\mathcal{N}(\mu, \sigma^2)$ -distributed and Y_1, \dots, Y_m be i.i.d. $\mathcal{N}(\nu, \tau^2)$ -distributed.

Case i If $\theta = (\mu, \nu, \sigma^2, \tau^2) \in \mathbb{R}^2 \times \mathbb{R}_+^2$, one can easily check that

$$S := \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2, \sum_{j=1}^m Y_j, \sum_{j=1}^m Y_j^2 \right)$$

is sufficient and complete.

Case ii If μ, σ^2 and τ^2 are unknown, and $\nu = \mu$, then S of course remains sufficient. One can however show that S is not complete. Difficult question: does a complete statistic exist?

3.3 The Cramer-Rao lower bound

Let $\{P_\theta : \theta \in \Theta\}$ be a collection of distributions on \mathcal{X} , dominated by a σ -finite measure ν . We denote the densities by

$$p_\theta := \frac{dP_\theta}{d\nu}, \quad \theta \in \Theta.$$

In this section, we assume that Θ is a one-dimensional open interval (the extension to a higher-dimensional parameter space will be handled in the next section).

We will impose the following two conditions:

Condition I *The set*

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

does not depend on θ .

Condition II (*Differentiability in L_2*) *For all θ and for a function $s_\theta : \mathcal{X} \rightarrow \mathbb{R}$ satisfying*

$$I(\theta) := E_\theta s_\theta(X)^2 < \infty,$$

it holds that

$$\lim_{h \rightarrow 0} E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{hp_\theta(X)} - s_\theta(X) \right)^2 = 0.$$

Definition *If I and II hold, we call s_θ the score function, and $I(\theta)$ the Fisher information.*

Lemma 3.3.1 *Assume conditions I and II. Then*

$$E_\theta s_\theta(X) = 0, \quad \forall \theta.$$

Proof. Under P_θ , we only need to consider values x with $p_\theta(x) > 0$, that is, we may freely divide by p_θ , without worrying about dividing by zero.

Observe that

$$E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{p_\theta(X)} \right) = \int_A (p_{\theta+h} - p_\theta) d\nu = 0,$$

since densities integrate to 1, and both $p_{\theta+h}$ and p_θ vanish outside A . Thus,

$$\begin{aligned} |E_\theta s_\theta(X)|^2 &= \left| E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{hp_\theta(X)} - s_\theta(X) \right) \right|^2 \\ &\leq E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{hp_\theta(X)} - s_\theta(X) \right)^2 \rightarrow 0. \end{aligned}$$

□

Note Thus $I(\theta) = \text{var}_\theta(s_\theta(X))$.

Remark If $p_\theta(x)$ is differentiable for all x , we take

$$s_\theta(x) := \frac{d}{d\theta} \log p_\theta(x) = \frac{\dot{p}_\theta(x)}{p_\theta(x)},$$

where

$$\dot{p}_\theta(x) := \frac{d}{d\theta} p_\theta(x).$$

Remark Suppose X_1, \dots, X_n are i.i.d. with density p_θ , and $s_\theta = \dot{p}_\theta/p_\theta$ exists. The joint density is

$$\mathbf{p}_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i),$$

so that (under conditions I and II) the score function for n observations is

$$\mathbf{s}_\theta(\mathbf{x}) = \sum_{i=1}^n s_\theta(x_i).$$

The Fisher information for n observations is thus

$$\mathbf{I}(\theta) = \text{var}_\theta(\mathbf{s}_\theta(\mathbf{X})) = \sum_{i=1}^n \text{var}_\theta(s_\theta(X_i)) = nI(\theta).$$

Theorem 3.3.1 (*The Cramer-Rao lower bound*) Suppose conditions I and II are met, and that T is an unbiased estimator of $g(\theta)$ with finite variance. Then $g(\theta)$ has a derivative, $\dot{g}(\theta) := dg(\theta)/d\theta$, equal to

$$\dot{g}(\theta) = \text{cov}(T, s_\theta(X)).$$

Moreover,

$$\text{var}_\theta(T) \geq \frac{[\dot{g}(\theta)]^2}{I(\theta)}, \quad \forall \theta.$$

Proof. We first show differentiability of g . As T is unbiased, we have

$$\begin{aligned} \frac{g(\theta+h) - g(\theta)}{h} &= \frac{E_{\theta+h}T(X) - E_{\theta}T(X)}{h} \\ &= \frac{1}{h} \int T(p_{\theta+h} - p_{\theta})d\nu = E_{\theta}T(X) \frac{p_{\theta+h}(X) - p_{\theta}(X)}{hp_{\theta}(X)} \\ &= E_{\theta}T(X) \left(\frac{p_{\theta+h}(X) - p_{\theta}(X)}{hp_{\theta}(X)} - s_{\theta}(X) \right) + E_{\theta}T(X)s_{\theta}(X) \\ &= E_{\theta} \left(T(X) - g_{\theta} \right) \left(\frac{p_{\theta+h}(X) - p_{\theta}(X)}{hp_{\theta}(X)} - s_{\theta}(X) \right) + E_{\theta}T(X)s_{\theta}(X) \\ &\quad \rightarrow E_{\theta}T(X)s_{\theta}(X), \end{aligned}$$

as, by the Cauchy-Schwarz inequality

$$\begin{aligned} &\left| E_{\theta} \left(T(X) - g_{\theta} \right) \left(\frac{p_{\theta+h}(X) - p_{\theta}(X)}{hp_{\theta}(X)} - s_{\theta}(X) \right) \right|^2 \\ &\leq \text{var}_{\theta}(T) E_{\theta} \left(\frac{p_{\theta+h}(X) - p_{\theta}(X)}{hp_{\theta}(X)} - s_{\theta}(X) \right)^2 \rightarrow 0. \end{aligned}$$

Thus,

$$\dot{g}(\theta) = E_{\theta}T(X)s_{\theta}(X) = \text{cov}_{\theta}(T, s_{\theta}(X)).$$

The last inequality holds because $E_{\theta}s_{\theta}(X) = 0$. By Cauchy-Schwarz,

$$\begin{aligned} [\dot{g}(\theta)]^2 &= |\text{cov}_{\theta}(T, s_{\theta}(X))|^2 \\ &\leq \text{var}_{\theta}(T)\text{var}_{\theta}(s_{\theta}(X)) = \text{var}_{\theta}(T)I(\theta). \end{aligned}$$

□

Definition We call $[\dot{g}(\theta)]^2/I(\theta)$, $\theta \in \Theta$, the Cramer Rao lower bound (CRLB) (for estimating $g(\theta)$).

Example 3.3.1 Let X_1, \dots, X_n be i.i.d. Exponential(θ), $\theta > 0$. The density of a single observation is then

$$p_{\theta}(x) = \theta e^{-\theta x}, \quad x > 0.$$

Let $g(\theta) := 1/\theta$, and $T := \bar{X}$. Then T is unbiased, and $\text{var}_{\theta}(T) = 1/(n\theta^2)$. We now compute the CRLB. With $g(\theta) = 1/\theta$, one has $\dot{g}(\theta) = -1/\theta^2$. Moreover,

$$\log p_{\theta}(x) = \log \theta - \theta x,$$

so

$$s_{\theta}(x) = 1/\theta - x,$$

and hence

$$I(\theta) = \text{var}_{\theta}(X) = \frac{1}{\theta^2}.$$

The CRLB for n observations is thus

$$\frac{[\dot{g}(\theta)]^2}{nI(\theta)} = \frac{1}{n\theta^2}.$$

In other words, T reaches the CRLB.

Example 3.3.2 Suppose X_1, \dots, X_n are i.i.d. $\text{Poisson}(\theta)$, $\theta > 0$. Then

$$\log p_\theta(x) = -\theta + x \log \theta - \log x!,$$

so

$$s_\theta(x) = -1 + \frac{x}{\theta},$$

and hence

$$I(\theta) = \text{var}_\theta \left(\frac{X}{\theta} \right) = \frac{\text{var}_\theta(X)}{\theta^2} = \frac{1}{\theta}.$$

One easily checks that \bar{X} reaches the CRLB for estimating θ .

Let now $g(\theta) := e^{-\theta}$. The UMVU estimator of $g(\theta)$ is

$$T := \left(1 - \frac{1}{n} \right)^{\sum_{i=1}^n X_i}.$$

To compute its variance, we first compute

$$\begin{aligned} E_\theta T^2 &= \sum_{k=0}^{\infty} \left(1 - \frac{1}{n} \right)^{2k} \frac{(n\theta)^k}{k!} e^{-n\theta} \\ &= e^{-n\theta} \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{(n-1)^2 \theta}{n} \right)^k \\ &= e^{-n\theta} \exp \left[\frac{(n-1)^2 \theta}{n} \right] = \exp \left[\frac{(1-2n)\theta}{n} \right]. \end{aligned}$$

Thus,

$$\begin{aligned} \text{var}_\theta(T) &= E_\theta T^2 - [E_\theta T]^2 = E_\theta T^2 - e^{-2\theta} \\ &= e^{-2\theta} \left(e^{\theta/n} - 1 \right) \\ &\begin{cases} > \theta e^{-2\theta}/n \\ \approx \theta e^{-2\theta}/n \text{ for } n \text{ large} \end{cases}. \end{aligned}$$

As $\dot{g}(\theta) = -e^{-\theta}$, the CRLB is

$$\frac{[\dot{g}(\theta)]^2}{nI(\theta)} = \frac{\theta e^{-2\theta}}{n}.$$

We conclude that T does not reach the CRLB, but the gap is small for n large.

For the next result, we:

Recall Let X and Y be two real-valued random variables. The correlation between X and Y is

$$\rho(X, Y) := \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}}.$$

We have

$$|\rho(X, Y)| = 1 \Leftrightarrow \exists \text{ constants } a, b : Y = aX + b \text{ (a.s.)}.$$

The next lemma shows that the CRLB can only be reached within exponential families, thus is only tight in a rather limited context.

Lemma 3.3.2 *Assume conditions I and II, with $s_\theta = \dot{p}_\theta/p_\theta$. Suppose T is unbiased for $g(\theta)$, and that T reaches the Cramer-Rao lower bound. Then $\{P_\theta : \theta \in \Theta\}$ forms a one-dimensional exponential family: there exist functions $c(\theta)$, $d(\theta)$, and $h(x)$ such that for all θ ,*

$$p_\theta(x) = \exp[c(\theta)T(X) - d(\theta)]h(x), \quad x \in \mathcal{X}.$$

Moreover, $c(\theta)$ and $d(\theta)$ are differentiable, say with derivatives $\dot{c}(\theta)$ and $\dot{d}(\theta)$ respectively. We furthermore have the equality

$$g(\theta) = \dot{d}(\theta)/\dot{c}(\theta), \quad \forall \theta.$$

Proof. By Theorem 3.3.1, when T reaches the CRLB, we must have

$$\text{var}_\theta(T) = \frac{|\text{cov}(T, s_\theta(X))|^2}{\text{var}_\theta(s_\theta(X))},$$

i.e., then the correlation between T and $s_\theta(X)$ is ± 1 . Thus, there exist constants $a(\theta)$ and $b(\theta)$ (depending on θ), such that

$$s_\theta(X) = a(\theta)T(X) - b(\theta). \quad (3.1)$$



But, as $s_\theta = \dot{p}_\theta/p_\theta = d \log p_\theta / d\theta$, we can take primitives:

$$\log p_\theta(x) = c(\theta)T(x) - d(\theta) + \tilde{h}(x),$$

where $\dot{c}(\theta) = a(\theta)$, $\dot{d}(\theta) = b(\theta)$ and $\tilde{h}(x)$ is constant in θ . Hence,

$$p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x),$$

with $h(x) = \exp[\tilde{h}(x)]$.

Moreover, the equation (3.1) tells us that

$$E_\theta s_\theta(X) = a(\theta)E_\theta T - b(\theta) = a(\theta)g(\theta) - b(\theta).$$

Because $E_\theta s_\theta(X) = 0$, this implies that $g(\theta) = b(\theta)/a(\theta)$. □

3.4 Higher-dimensional extensions

Expectations and covariances of random vectors

Let $X \in \mathbb{R}^p$ be a p -dimensional random vector. Then EX is a p -dimensional vector, and

$$\Sigma := \text{Cov}(X) := EXX^T - (EX)(EX)^T$$

is a $p \times p$ matrix containing all variances (on the diagonal) and covariances (off-diagonal). Note that Σ is positive semi-definite: for any vector $a \in \mathbb{R}^p$, we have

$$\text{var}(a^T X) = a^T \Sigma a \geq 0.$$

Some matrix algebra

Let V be a symmetric matrix. If V is positive (semi-)definite, we write this as $V > 0$ ($V \geq 0$). One then has that $V = W^2$, where W is also positive (semi-)definite.

Auxiliary lemma. *Suppose $V > 0$. Then*

$$\max_{a \in \mathbb{R}^p} \frac{|a^T c|^2}{a^T V a} = c^T V^{-1} c.$$

Proof. Write $V = W^2$, and $b := W a$, $d := W^{-1} c$. Then $a^T V a = b^T b = \|b\|^2$ and $a^T c = b^T d$. By Cauchy-Schwarz

$$\max_{b \in \mathbb{R}^p} \frac{|b^T d|^2}{\|b\|^2} = \|d\|^2 = d^T d = c^T V^{-1} c.$$

□

We will now present the CRLB in higher dimensions. To simplify the exposition, we will not carefully formulate the regularity conditions, that is, we assume derivatives to exist and that we can interchange differentiation and integration at suitable places.



Consider a parameter space $\Theta \subset \mathbb{R}^p$. Let

$$g : \Theta \rightarrow \mathbb{R},$$

be a given function. Denote the vector of partial derivatives as

$$\dot{g}(\theta) := \begin{pmatrix} \partial g(\theta) / \partial \theta_1 \\ \vdots \\ \partial g(\theta) / \partial \theta_p \end{pmatrix}.$$

The score vector is defined as

$$s_\theta(\cdot) := \begin{pmatrix} \partial \log p_\theta / \partial \theta_1 \\ \vdots \\ \partial \log p_\theta / \partial \theta_p \end{pmatrix}.$$

The Fisher information matrix is

$$I(\theta) = E_\theta s_\theta(X) s_\theta^T(X) = \text{Cov}_\theta(s_\theta(X)).$$

Theorem 3.4.1 *Let T be an unbiased estimator of $g(\theta)$. Then, under regularity conditions,*

$$\text{var}_\theta(T) \geq \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta).$$



Proof. As in the one-dimensional case, one can show that, for $j = 1, \dots, p$,

$$\dot{g}_j(\theta) = \text{cov}_\theta(T, s_{\theta,j}(X)).$$

Hence, for all $a \in \mathbb{R}^p$,

$$\begin{aligned} |a^T \dot{g}(\theta)|^2 &= |\text{cov}_\theta(T, a^T s_\theta(X))|^2 \\ &\leq \text{var}_\theta(T) \text{var}_\theta(a^T s_\theta(X)) \\ &= \text{var}_\theta(T) a^T I(\theta) a. \end{aligned}$$

Combining this with the auxiliary lemma gives

$$\text{var}_\theta(T) \geq \max_{a \in \mathbb{R}^p} \frac{|a^T \dot{g}(\theta)|^2}{a^T I(\theta) a} = \dot{g}(\theta)^T I(\theta)^{-1} \dot{g}(\theta).$$

□

Corollary 3.4.1 *As a consequence, one obtains a lower bound for unbiased estimators of higher-dimensional parameters of interest. As example, let $g(\theta) := \theta = (\theta_1, \dots, \theta_p)^T$, and suppose that $T \in \mathbb{R}^p$ is an unbiased estimator of θ . Then, for all $a \in \mathbb{R}^p$, $a^T T$ is an unbiased estimator of $a^T \theta$. Since $a^T \theta$ has derivative a , the CRLB gives*

$$\text{var}_\theta(a^T T) \geq a^T I(\theta)^{-1} a.$$

But

$$\text{var}_\theta(a^T T) = a^T \text{Cov}_\theta(T) a.$$

So for all a ,

$$a^T \text{Cov}_\theta(T) a \geq a^T I(\theta)^{-1} a,$$

in other words, $\text{Cov}_\theta(T) \geq I(\theta)^{-1}$, that is, $\text{Cov}_\theta(T) - I(\theta)^{-1}$ is positive semi-definite.

3.5 Uniformly most powerful tests

3.5.1 An example

Let X_1, \dots, X_n be i.i.d. copies of a Bernoulli random variable $X \in \{0, 1\}$ with success parameter $\theta \in (0, 1)$:

$$P_\theta(X = 1) = 1 - P_\theta(X = 0) = \theta.$$

We consider three testing problems. The chosen level in all three problems is $\alpha = 0.05$.

Problem 1

We want to test, at level α , the hypothesis

$$H_0 : \theta = 1/2 := \theta_0,$$

against the alternative

$$H_1 : \theta = 1/4 := \theta_1.$$

Let $T := \sum_{i=1}^n X_i$ be the number of successes (T is a sufficient statistic), and consider the randomized test

$$\phi(T) := \begin{cases} 1 & \text{if } T < t_0 \\ q & \text{if } T = t_0 \\ 0 & \text{if } T > t_0 \end{cases},$$

where $q \in (0, 1)$, and where t_0 is the critical value of the test. The constants q and $t_0 \in \{0, \dots, n\}$ are chosen in such a way that the probability of rejecting H_0 when it is in fact true, is equal to α :

$$P_{\theta_0}(H_0 \text{ rejected}) = P_{\theta_0}(T \leq t_0 - 1) + qP_{\theta_0}(T = t_0) := \alpha.$$

Thus, we take t_0 in such a way that

$$P_{\theta_0}(T \leq t_0 - 1) \leq \alpha, \quad P_{\theta_0}(T \leq t_0) > \alpha,$$

(i.e., $t_0 - 1 = q_+(\alpha)$ with q_+ the quantile function defined in Section 1.6) and

$$q = \frac{\alpha - P_{\theta_0}(T \leq t_0 - 1)}{P_{\theta_0}(T = t_0)}.$$

Because $\phi = \phi_{\text{NP}}$ is the Neyman Pearson test, it is the most powerful test (at level α) (see the Neyman Pearson Lemma in Section 2.2). The power of the test is $\beta(\theta_1)$, where

$$\beta(\theta) := E_{\theta}\phi(T).$$

Numerical Example

Let $n = 7$. Then

$$P_{\theta_0}(T = 0) = \left(\frac{1}{2}\right)^7 = 0.0078,$$

$$P_{\theta_0}(T = 1) = \binom{7}{1} \left(\frac{1}{2}\right)^7 = 0.0546,$$

$$P_{\theta_0}(T \leq 1) = 0.0624 > \alpha,$$

so we choose $t_0 = 1$. Moreover

$$q = \frac{0.05 - 0.0078}{0.0546} = \frac{422}{546}.$$

The power is now

$$\begin{aligned}\beta(\theta_1) &= P_{\theta_1}(T = 0) + qP_{\theta_1}(T = 1) \\ &= \left(\frac{3}{4}\right)^7 + \frac{422}{546} \binom{7}{1} \left(\frac{3}{4}\right)^6 \left(\frac{1}{4}\right) = 0.1335 + \frac{422}{546} 0.3114.\end{aligned}$$

Problem 2

Consider now testing

$$H_0 : \theta_0 = 1/2,$$

against

$$H_1 : \theta < 1/2.$$

In Problem 1, the construction of the test ϕ is independent of the value $\theta_1 < \theta_0$. So ϕ is most powerful for all $\theta_1 < \theta_0$. We say that ϕ is *uniformly most powerful* (German: *gleichmässig mächtigst*) for the alternative $H_1 : \theta < \theta_0$.

Problem 3

We now want to test

$$H_0 : \theta \geq 1/2,$$

against the alternative

$$H_1 : \theta < 1/2.$$

Recall the function

$$\beta(\theta) := E_{\theta}\phi(T).$$

The level of ϕ is defined as

$$\sup_{\theta \geq 1/2} \beta(\theta).$$

We have

$$\begin{aligned}\beta(\theta) &= P_{\theta}(T \leq t_0 - 1) + qP_{\theta}(T = t_0) \\ &= (1 - q)P_{\theta}(T \leq t_0 - 1) + qP_{\theta}(T \leq t_0).\end{aligned}$$

Observe that if $\theta_1 < \theta_0$, small values of T are more likely under P_{θ_1} than under P_{θ_0} :

$$P_{\theta_1}(T \leq t) > P_{\theta_0}(T \leq t), \quad \forall t \in \{0, 1, \dots, n\}.$$

Thus, $\beta(\theta)$ is a decreasing function of θ . It follows that the level of ϕ is

$$\sup_{\theta \geq 1/2} \beta(\theta) = \beta(1/2) = \alpha.$$

Hence, ϕ is uniformly most powerful for $H_0 : \theta \geq 1/2$ against $H_1 : \theta < 1/2$.

3.5.2 UMP tests and exponential families

Let $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ be a family of probability measures. Let $\Theta_0 \subset \Theta$, $\Theta_1 \subset \Theta$, and $\Theta_0 \cap \Theta_1 = \emptyset$. Based on observations X , with distribution $P \in \mathcal{P}$, we consider the general testing problem, at level α , for

$$H_0 : \theta \in \Theta_0,$$

against

$$H_1 : \theta \in \Theta_1.$$

We say that a test ϕ has level α if

$$\sup_{\theta \in \Theta_0} E_\theta \phi(X) \leq \alpha.$$

Definition A test ϕ is called Uniformly Most Powerful (UMP) if

- ϕ has level α ,
- for all tests ϕ' with level α , it holds that $E_\theta \phi'(X) \leq E_\theta \phi(X) \forall \theta \in \Theta_1$.

We now simplify the situation to the case where Θ is an interval in \mathbb{R} , and to the testing problem

$$H_0 : \theta \leq \theta_0,$$

against

$$H_1 : \theta > \theta_0.$$

We also suppose that \mathcal{P} is dominated by a σ -finite measure ν .

Theorem 3.5.1 Suppose that \mathcal{P} is a one-dimensional exponential family

$$\frac{dP_\theta}{d\nu}(x) := p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x).$$

Assume moreover that $c(\theta)$ is a strictly increasing function of θ . Then the UMP test ϕ is

$$\phi(T(x)) := \begin{cases} 1 & \text{if } T(x) > t_0 \\ q & \text{if } T(x) = t_0 \\ 0 & \text{if } T(x) < t_0 \end{cases},$$

where q and t_0 are chosen in such a way that $E_{\theta_0} \phi(T) = \alpha$.

Proof. The Neyman Pearson test for $H_0 : \theta = \theta_0$ against $H_1 : \theta = \theta_1$ is

$$\phi_{\text{NP}}(x) = \begin{cases} 1 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) > c_0 \\ q_0 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) = c_0 \\ 0 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) < c_0 \end{cases},$$

where q_0 and c_0 are chosen in such a way that $E_{\theta_0} \phi_{\text{NP}}(X) = \alpha$. We have

$$\frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} = \exp \left[(c(\theta_1) - c(\theta_0))T(X) - (d(\theta_1) - d(\theta_0)) \right].$$

Hence

$$\frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} \begin{matrix} > \\ = c \\ < \end{matrix} \Leftrightarrow T(x) = t \begin{matrix} > \\ \\ < \end{matrix},$$

where t is some constant (depending on c , θ_0 and θ_1). Therefore, $\phi = \phi_{\text{NP}}$. It follows that ϕ is most powerful for $H_0 : \theta = \theta_0$ against $H_1 : \theta = \theta_1$. Because ϕ does not depend on θ_1 , it is therefore UMP for $H_0 : \theta = \theta_0$ against $H_1 : \theta > \theta_0$.

We will now prove that $\beta(\theta) := E_{\theta}\phi(T)$ is increasing in θ . Let

$$\bar{p}_{\theta}(t) = \exp[c(\theta)t - d(\theta)],$$

This is the density of T with respect to

$$d\bar{\nu}(t) = \int_{x: T(x)=t} h(x)d\nu(x).$$

For $\vartheta > \theta$

$$\frac{\bar{p}_{\vartheta}(t)}{\bar{p}_{\theta}(t)} = \exp\left[(c(\vartheta) - c(\theta))t - (d(\vartheta) - d(\theta))\right],$$

which is increasing in t . Moreover, we have

$$\int \bar{p}_{\vartheta}d\bar{\nu} = \int \bar{p}_{\theta}d\bar{\nu} = 1.$$

Therefore, there must be a point s_0 where the two densities cross:

$$\begin{cases} \frac{\bar{p}_{\vartheta}(t)}{\bar{p}_{\theta}(t)} \leq 1 & \text{for } t \leq s_0 \\ \frac{\bar{p}_{\vartheta}(t)}{\bar{p}_{\theta}(t)} \geq 1 & \text{for } t \geq s_0 \end{cases}.$$

But then

$$\begin{aligned} \beta(\vartheta) - \beta(\theta) &= \int \phi(t)[\bar{p}_{\vartheta}(t) - \bar{p}_{\theta}(t)]d\bar{\nu}(t) \\ &= \int_{t \leq s_0} \phi(t)[\bar{p}_{\vartheta}(t) - \bar{p}_{\theta}(t)]d\bar{\nu}(t) + \int_{t \geq s_0} \phi(t)[\bar{p}_{\vartheta}(t) - \bar{p}_{\theta}(t)]d\bar{\nu}(t) \\ &\geq \phi(s_0) \int [\bar{p}_{\vartheta}(t) - \bar{p}_{\theta}(t)]d\bar{\nu}(t) = 0. \end{aligned}$$

So indeed $\beta(\theta)$ is increasing in θ .

But then

$$\sup_{\theta \leq \theta_0} \beta(\theta) = \beta(\theta_0) = \alpha.$$

Hence, ϕ has level α . Because any other test ϕ' with level α must have $E_{\theta_0}\phi'(X) \leq \alpha$, we conclude that ϕ is UMP.

□

Example 3.5.1 Let X_1, \dots, X_n be an i.i.d. sample from the $\mathcal{N}(\mu_0, \sigma^2)$ -distribution, with μ_0 known, and $\sigma^2 > 0$ unknown. We want to test

$$H_0 : \sigma^2 \leq \sigma_0^2,$$

against

$$H_1 : \sigma^2 > \sigma_0^2.$$

The density of the sample is

$$\mathbf{p}_{\sigma^2}(x_1, \dots, x_n) = \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2 - \frac{n}{2} \log(2\pi\sigma^2) \right].$$

Thus, we may take

$$c(\sigma^2) = -\frac{1}{2\sigma^2},$$

and

$$T(\mathbf{X}) = \sum_{i=1}^n (X_i - \mu_0)^2.$$

The function $c(\sigma^2)$ is strictly increasing in σ^2 . So we let ϕ be the test which rejects H_0 for large values of $T(\mathbf{X})$.

Example 3.5.2 Let X_1, \dots, X_n be an i.i.d. sample from the Bernoulli(θ)-distribution, $0 < \theta < 1$. Then

$$\mathbf{p}_{\theta}(x_1, \dots, x_n) = \exp \left[\log \left(\frac{\theta}{1-\theta} \right) \sum_{i=1}^n x_i + \log(1-\theta) \right].$$

We can take

$$c(\theta) = \log \left(\frac{\theta}{1-\theta} \right),$$

which is strictly increasing in θ . Then $T(\mathbf{X}) = \sum_{i=1}^n X_i$.

Right-sided alternative

$$H_0 : \theta \leq \theta_0 ,$$

against

$$H_1 : \theta > \theta_0 .$$

The UMP test is

$$\phi_R(T) := \begin{cases} 1 & T > t_R \\ q_R & T = t_R \\ 0 & T < t_R \end{cases} .$$

The function $\beta_R(\theta) := E_{\theta}\phi_R(T)$ is strictly increasing in θ .

Left-sided alternative

$$H_0 : \theta \geq \theta_0 ,$$

against

$$H_1 : \theta < \theta_0 .$$

The UMP test is

$$\phi_L(T) := \begin{cases} 1 & T < t_L \\ q_L & T = t_L \\ 0 & T > t_L \end{cases} .$$

The function $\beta_L(\theta) := E_\theta \phi_L(T)$ is strictly decreasing in θ .

Two-sided alternative

$$H_0 : \theta = \theta_0 ,$$

against

$$H_1 : \theta \neq \theta_0 .$$

The test ϕ_R is most powerful for $\theta > \theta_0$, whereas ϕ_L is most powerful for $\theta < \theta_0$. Hence, a UMP test does not exist for the two-sided alternative.

3.5.3 Unbiased tests

Consider again the general case: $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ is a family of probability measures, the spaces Θ_0 , and Θ_1 are disjoint subspaces of Θ , and the testing problem is

$$H_0 : \theta \in \Theta_0 ,$$

against

$$H_1 : \theta \in \Theta_1 .$$

The significance level is α ($\alpha < 1$).

As we have seen in Example 3.5.2, uniformly most powerful tests do not always exist. We therefore restrict attention to a smaller class of tests, and look for uniformly most powerful tests in the smaller class.

Definition A test ϕ is called unbiased (German unverfälscht) if for all $\theta \in \Theta_0$ and all $\vartheta \in \Theta_1$,

$$E_\theta \phi(X) \leq E_\vartheta \phi(X) .$$

Definition A test ϕ is called Uniformly Most Powerful Unbiased (UMPU) if

- ϕ has level α ,
- ϕ is unbiased,
- for all unbiased tests ϕ' with level α , one has $E_\theta \phi'(X) \leq E_\theta \phi(X) \forall \theta \in \Theta_1$.

We return to the special case where $\Theta \subset \mathbb{R}$ is an interval. We consider testing

$$H_0 : \theta = \theta_0 ,$$

against

$$H_1 : \theta \neq \theta_0 .$$

The following theorem presents the UMPU test. We omit the proof (see e.g. Lehmann ...).

Theorem 3.5.2 *Suppose \mathcal{P} is a one-dimensional exponential family:*

$$\frac{dP_\theta}{d\nu}(x) := p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)]h(x),$$

with $c(\theta)$ strictly increasing in θ . Then the UMPU test is

$$\phi(T(x)) := \begin{cases} 1 & \text{if } T(x) < t_L \text{ or } T(x) > t_R \\ q_L & \text{if } T(x) = t_L \\ q_R & \text{if } T(x) = t_R \\ 0 & \text{if } t_L < T(x) < t_R \end{cases},$$

where the constants t_R , t_L , q_R and q_L are chosen in such a way that

$$E_{\theta_0}\phi(X) = \alpha, \quad \left. \frac{d}{d\theta} E_\theta\phi(X) \right|_{\theta=\theta_0} = 0.$$

Note Let ϕ_R a right-sided test as defined Theorem 3.5.1 with level at most α and ϕ_L be the similarly defined left-sided test. Then $\beta_R(\theta) = E_\theta\phi_R(T)$ is strictly increasing, and $\beta_L(\theta) = E_\theta\phi_L(T)$ is strictly decreasing. The two-sided test ϕ of Theorem 3.5.2 is a superposition of two one-sided tests. Writing

$$\beta(\theta) = E_\theta\phi(T),$$

the one-sided tests are constructed in such a way that

$$\beta(\theta) = \beta_R(\theta) + \beta_L(\theta).$$

Moreover, $\beta(\theta)$ should be minimal at $\theta = \theta_0$, whence the requirement that its derivative at θ_0 should vanish. Let us see what this derivative looks like. With the notation used in the proof of Theorem 3.5.1, for a test $\tilde{\phi}$ depending only on the sufficient statistic T ,

$$E_\theta\tilde{\phi}(T) = \int \tilde{\phi}(t) \exp[c(\theta)t - d(\theta)]d\bar{\nu}(t).$$

Hence, assuming we can take the differentiation inside the integral,

$$\begin{aligned} \frac{d}{d\theta} E_\theta\tilde{\phi}(T) &= \int \tilde{\phi}(t) \exp[c(\theta)t - d(\theta)](\dot{c}(\theta)t - \dot{d}(\theta))d\bar{\nu}(t) \\ &= \dot{c}(\theta)\text{cov}_\theta(\tilde{\phi}(T), T). \end{aligned}$$

Example 3.5.3 Let X_1, \dots, X_n be an i.i.d. sample from the $\mathcal{N}(\mu, \sigma_0^2)$ -distribution, with $\mu \in \mathbb{R}$ unknown, and with σ_0^2 known. We consider testing

$$H_0 : \mu = \mu_0,$$

against

$$H_1 : \mu \neq \mu_0.$$

A sufficient statistic is $T := \sum_{i=1}^n X_i$. We have, for $t_L < t_R$,

$$\begin{aligned} E_\mu \phi(T) &= \mathbb{P}_\mu(T > t_R) + \mathbb{P}_\mu(T < t_L) \\ &= \mathbb{P}_\mu\left(\frac{T - n\mu}{\sqrt{n}\sigma_0} > \frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) + \mathbb{P}_\mu\left(\frac{T - n\mu}{\sqrt{n}\sigma_0} < \frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right) \\ &= 1 - \Phi\left(\frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) + \Phi\left(\frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right), \end{aligned}$$

where Φ is the standard normal distribution function. To avoid confusion with the test ϕ , we denote the standard normal density in this example by $\dot{\Phi}$. Thus,

$$\frac{d}{d\mu} E_\mu \phi(T) = \frac{1}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) - \frac{1}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right),$$

So putting

$$\left. \frac{d}{d\mu} E_\mu \phi(T) \right|_{\mu=\mu_0} = 0,$$

gives

$$\dot{\Phi}\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) = \dot{\Phi}\left(\frac{t_L - n\mu_0}{\sqrt{n}\sigma_0}\right),$$

or

$$(t_R - n\mu_0)^2 = (t_L - n\mu_0)^2.$$

We take the solution $(t_L - n\mu_0) = -(t_R - n\mu_0)$, (because the solution $(t_L - n\mu_0) = (t_R - n\mu_0)$ leads to a test that always rejects, and hence does not have level α , as $\alpha < 1$). Plugging this solution back in gives

$$\begin{aligned} E_{\mu_0} \phi(T) &= 1 - \Phi\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) + \Phi\left(-\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) \\ &= 2 \left(1 - \Phi\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right)\right). \end{aligned}$$

The requirement $E_{\mu_0} \phi(T) = \alpha$ gives us

$$\Phi\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) = 1 - \alpha/2,$$

and hence

$$t_R - n\mu_0 = \sqrt{n}\sigma_0 \Phi^{-1}(1 - \alpha/2), \quad t_L - n\mu_0 = -\sqrt{n}\sigma_0 \Phi^{-1}(1 - \alpha/2).$$

3.5.4 Conditional tests

We now study the case where Θ is an interval in \mathbb{R}^2 . We let $\theta = (\beta, \gamma)$, and we assume that γ is the parameter of interest. We aim at testing

$$H_0 : \gamma \leq \gamma_0,$$

against the alternative

$$H_1 : \gamma > \gamma_0.$$

We assume moreover that we are dealing with an exponential family in canonical form:

$$p_\theta(x) = \exp[\beta T_1(x) + \gamma T_2(x) - d(\theta)]h(x).$$

Then we can restrict ourselves to tests $\phi(T)$ depending only on the sufficient statistic $T = (T_1, T_2)$.

Lemma 3.5.1 *Suppose that $\{\beta : (\beta, \gamma_0) \in \Theta\}$ contains an open interval. Let*

$$\phi(T_1, T_2) = \begin{cases} 1 & \text{if } T_2 > t_0(T_1) \\ q(T_1) & \text{if } T_2 = t_0(T_1) \\ 0 & \text{if } T_2 < t_0(T_1) \end{cases},$$

where the constants $t_0(T_1)$ and $q(T_1)$ are allowed to depend on T_1 , and are chosen in such a way that

$$E_{\gamma_0} \left(\phi(T_1, T_2) \middle| T_1 \right) = \alpha.$$

Then ϕ is UMPU.

Proof.



Let $\bar{p}_\theta(t_1, t_2)$ be the density of (T_1, T_2) with respect to $\bar{\nu}$:

$$\bar{p}_\theta(t_1, t_2) := \exp[\beta t_1 + \gamma t_2 - d(\theta)],$$

$$d\bar{\nu}(t_1, t_2) := \int_{T_1(x)=t_1, T_2(x)=t_2} h(x) d\nu(x).$$

The conditional density of T_2 given $T_1 = t_1$ is

$$\begin{aligned} \bar{p}_\theta(t_2|t_1) &= \frac{\exp[\beta t_1 + \gamma t_2 - d(\theta)]}{\int_{s_2} \exp[\beta t_1 + \gamma s_2 - d(\theta)] d\bar{\nu}(t_1, s_2)} \\ &= \exp[\gamma t_2 - d(\gamma|t_1)], \end{aligned}$$

where

$$d(\gamma|t_1) := \log \left(\int_{s_2} \exp[\gamma s_2] d\bar{\nu}(t_1, s_2) \right).$$

In other words, the conditional distribution of T_2 given $T_1 = t_1$ - does not depend on β ,

- is a one-parameter exponential family in canonical form.
This implies that given $T_1 = t_1$, ϕ is UMPU.

Result 1 *The test ϕ has level α , i.e.*

$$\sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi(T) = E_{(\beta, \gamma_0)} \phi(T) = \alpha, \quad \forall \beta.$$

Proof of Result 1.

$$\sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi(T) \geq E_{(\beta, \gamma_0)} \phi(T) = E_{(\beta, \gamma_0)} E_{\gamma_0}(\phi(T)|T_1) = \alpha.$$

Conversely,

$$\sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi(T) = \sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} E_{\gamma}(\phi(T)|T_1) \leq E_{(\beta, \gamma)} \sup_{\gamma \leq \gamma_0} E_{\gamma}(\phi(T)|T_1) = \alpha.$$

Result 2 *The test ϕ is unbiased.*

Proof of Result 2. If $\gamma > \gamma_0$, it holds that $E_{\gamma}(\phi(T)|T_1) \geq \alpha$, as the conditional test is unbiased. Thus, also, for all β ,

$$E_{(\beta, \gamma)} \phi(T) = E_{(\beta, \gamma)} E_{\gamma}(\phi(T)|T_1) \geq \alpha,$$

i.e., ϕ is unbiased.

Result 3 *Let ϕ' be a test with level*

$$\alpha' := \sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi'(T) \leq \alpha,$$

and suppose moreover that ϕ' is unbiased, i.e., that

$$\sup_{\gamma \leq \gamma_0} \sup_{\beta} E_{(\beta, \gamma)} \phi'(T) \leq \inf_{\gamma > \gamma_0} \inf_{\beta} E_{(\beta, \gamma)} \phi'(T).$$

Then, conditionally on T_1 , ϕ' has level α' .

Proof of Result 3. As

$$\alpha' = \sup_{\gamma \leq \gamma_0} \sup_{\beta} E_{(\gamma, \beta)} \phi'(T)$$

we know that

$$E_{(\beta, \gamma_0)} \phi'(T) \leq \alpha', \quad \forall \beta.$$

Conversely, the unbiasedness implies that for all $\gamma > \gamma_0$,

$$E_{(\beta, \gamma)} \phi'(T) \geq \alpha', \quad \forall \beta.$$

A continuity argument therefore gives

$$E_{(\beta, \gamma_0)} \phi'(T) = \alpha', \quad \forall \beta.$$

In other words, we have

$$E_{(\beta, \gamma_0)}(\phi'(T) - \alpha') = 0, \forall \beta.$$

But then also

$$E_{(\beta, \gamma_0)} E_{\gamma_0} \left((\phi'(T) - \alpha') \middle| T_1 \right) = 0, \forall \beta,$$

which we can write as

$$E_{(\beta, \gamma_0)} h(T_1) = 0, \forall \beta.$$

The assumption that $\{\beta : (\beta, \gamma_0) \in \Theta\}$ contains an open interval implies that T_1 is complete for (β, γ_0) . So we must have

$$h(T_1) = 0, P_{(\beta, \gamma_0)}\text{-a.s.}, \forall \beta,$$

or, by the definition of h ,

$$E_{\gamma_0}(\phi'(T)|T_1) = \alpha', P_{(\beta, \gamma_0)}\text{-a.s.}, \forall \beta.$$

So conditionally on T_1 , the test ϕ' has level α' .

Result 4 *Let ϕ' be a test as given in Result 3. Then ϕ' can not be more powerful than ϕ at any (β, γ) , with $\gamma > \gamma_0$.*

Proof of Result 4. By the Neyman Pearson lemma, conditionally on T_1 , we have

$$E_{\gamma}(\phi'(T)|T_1) \geq E_{\gamma}(\phi(T)|T_1), \forall \gamma > \gamma_0.$$

Thus also

$$E_{(\beta, \gamma)} \phi'(T) \geq E_{(\beta, \gamma)} \phi(T), \forall \beta, \gamma > \gamma_0.$$

□

Example 3.5.4 Consider two independent samples $\mathbf{X} = (X_1, \dots, X_n)$ and $\mathbf{Y} = (Y_1, \dots, Y_m)$, where X_1, \dots, X_n are i.i.d. Poisson(λ)-distributed, and Y_1, \dots, Y_m are i.i.d. Poisson(μ)-distributed. We aim at testing

$$H_0 : \lambda \leq \mu,$$

against the alternative

$$H_1 : \lambda > \mu.$$

Define

$$\beta := \log(m\mu), \quad \gamma := \log((n\lambda)/(m\mu)).$$

The testing problem is equivalent to

$$H_0 : \gamma \leq \gamma_0,$$

against the alternative

$$H_1 : \gamma > \gamma_0,$$

where $\gamma_0 := \log(n/m)$.

The density is

$$\begin{aligned}
& \mathbf{p}_\theta(x_1, \dots, x_n, y_1, \dots, y_m) \\
&= \exp \left[\log(n\lambda) \sum_{i=1}^n x_i + \log(m\mu) \sum_{j=1}^m y_j - n\lambda - m\mu \right] \prod_{i=1}^n \frac{1}{x_i!} \prod_{j=1}^m \frac{1}{y_j!} \\
&= \exp \left[\log(m\mu) \left(\sum_{i=1}^n x_i + \sum_{j=1}^m y_j \right) + \log((n\lambda)/(m\mu)) \sum_{i=1}^n x_i - n\lambda - m\mu \right] h(\mathbf{x}, \mathbf{y}) \\
&= \exp[\beta T_1(\mathbf{x}, \mathbf{y}) + \gamma T_2(\mathbf{x}) - d(\theta)] h(\mathbf{x}, \mathbf{y}),
\end{aligned}$$

where

$$T_1(\mathbf{X}, \mathbf{Y}) := \sum_{i=1}^n X_i + \sum_{j=1}^m Y_j,$$

and

$$T_2(\mathbf{X}) := \sum_{i=1}^n X_i,$$

and

$$h(\mathbf{x}, \mathbf{y}) := \prod_{i=1}^n \frac{1}{x_i!} \prod_{j=1}^m \frac{1}{y_j!}.$$

The conditional distribution of T_2 given $T_1 = t_1$ is the Binomial(t_1, p)-distribution, with

$$p = \frac{n\lambda}{n\lambda + m\mu} = \frac{e^\gamma}{1 + e^\gamma}.$$

Thus, conditionally on $T_1 = t_1$, using the observation T_2 from the Binomial(t_1, p)-distribution, we test

$$H_0 : p \leq p_0,$$

against the alternative

$$H_1 : p > p_0,$$

where $p_0 := n/(n + m)$. This test is UMPU for the unconditional problem.

Chapter 4

Equivariant statistics

As we have seen in the previous chapter, it can be useful to restrict attention to a collection of statistics satisfying certain desirable properties. In Chapter 3, we restricted ourselves to unbiased estimators. In this chapter, equivariance will be the key concept.

The data consists of i.i.d. real-valued random variables X_1, \dots, X_n . We write $\mathbf{X} := (X_1, \dots, X_n)$. The density w.r.t. some dominating measure ν , of a single observation is denoted by p_θ . The density of \mathbf{X} is $\mathbf{p}_\theta(\mathbf{x}) = \prod_i p_\theta(x_i)$, $\mathbf{x} = (x_1, \dots, x_n)$.

Location model

Then $\theta \in \mathbb{R}$ is a location parameter, and we assume

$$X_i = \theta + \epsilon_i, \quad i = 1, \dots, n.$$

We are interested in estimating θ . Both the parameter space Θ , as well as the action space \mathcal{A} , are the real line \mathbb{R} .

Location-scale model

Here $\theta = (\mu, \sigma)$, with $\mu \in \mathbb{R}$ a location parameter and $\sigma > 0$ a scale parameter. We assume

$$X_i = \theta + \sigma \epsilon_i, \quad i = 1, \dots, n.$$

The parameter space Θ and action space \mathcal{A} are both $\mathbb{R} \times (0, \infty)$.

4.1 Equivariance in the location model

Definition A statistic $T = T(\mathbf{X})$ is called location equivariant if for all constants $c \in \mathbb{R}$ and all $\mathbf{x} = (x_1, \dots, x_n)$,

$$T(x_1 + c, \dots, x_n + c) = T(x_1, \dots, x_n) + c.$$

Examples

$$T = \begin{cases} \bar{X} \\ X_{(\frac{n+1}{2})} \\ \dots \end{cases} \quad (n \text{ odd}) .$$

Definition A loss function $L(\theta, a)$ is called location invariant if for all $c \in \mathbb{R}$,

$$L(\theta + c, a + c) = L(\theta, a), \quad (\theta, a) \in \mathbb{R}^2.$$

In this section we abbreviate location equivariance (invariance) to simply equivariance (invariance), and we assume throughout that the loss $L(\theta, a)$ is invariant.

Corollary If T is equivariant (and $L(\theta, a)$ is invariant), then

$$\begin{aligned} R(\theta, T) &= E_{\theta}L(\theta, T(\mathbf{X})) = E_{\theta}L(0, T(\mathbf{X}) - \theta) \\ &= E_{\theta}L(0, T(\mathbf{X} - \theta)) = E_{\theta}L_0[T(\varepsilon)], \end{aligned}$$

where $L_0[a] := L(0, a)$ and $\varepsilon := (\varepsilon_1, \dots, \varepsilon_n)$. Because the distribution of ε does not depend on θ , we conclude that the risk does not depend on θ . We may therefore omit the subscript θ in the last expression:

$$R(\theta, T) = EL_0[T(\varepsilon)].$$

Since for $\theta = 0$, we have the equality $\mathbf{X} = \varepsilon$ we may alternatively write

$$R(\theta, T) = E_0L_0[T(\mathbf{X})] = R(0, T).$$

Definition An equivariant statistic T is called uniform minimum risk equivariant (UMRE) if

$$R(\theta, T) = \min_{d \text{ equivariant}} R(\theta, d), \quad \forall \theta,$$

or equivalently,

$$R(0, T) = \min_{d \text{ equivariant}} R(0, d).$$

Lemma 4.1.1 Let $Y_i := X_i - X_n$, $i = 1, \dots, n$, and $\mathbf{Y} := (Y_1, \dots, Y_n)$. We have

$$T \text{ equivariant} \Leftrightarrow T = T(\mathbf{Y}) + X_n.$$

Proof.

(\Rightarrow) Trivial.

(\Leftarrow) Replacing \mathbf{X} by $\mathbf{X} + c$ leaves \mathbf{Y} unchanged (i.e. \mathbf{Y} is invariant). So $T(\mathbf{X} + c) = T(\mathbf{Y}) + X_n + c = T(\mathbf{X}) + c$. \square

Theorem 4.1.1 Let $Y_i := X_i - X_n$, $i = 1, \dots, n$, $\mathbf{Y} := (Y_1, \dots, Y_n)$, and define

$$T^*(\mathbf{Y}) := \arg \min_v E\left(L_0[v + \epsilon_n] \middle| \mathbf{Y}\right).$$

Moreover, let

$$T^*(\mathbf{X}) := T^*(\mathbf{Y}) + X_n.$$

Then T^* is UMRE.

Proof. First, note that the distribution of \mathbf{Y} does not depend on θ , so that T^* is indeed a statistic. It is also equivariant, by the previous lemma.

Let T be an equivariant statistic. Then $T(\mathbf{X}) = T(\mathbf{Y}) + X_n$. So

$$T(\mathbf{X}) - \theta = T(\mathbf{Y}) + \epsilon_n.$$

Hence

$$R(0, T) = EL_0[T(\mathbf{Y}) + \epsilon_n] = E\left[E\left(L_0[T(\mathbf{Y}) + \epsilon_n] \middle| \mathbf{Y}\right)\right].$$

But

$$E\left(L_0[T(\mathbf{Y}) + \epsilon_n] \middle| \mathbf{Y}\right) \geq \min_v E\left(L_0[v + \epsilon_n] \middle| \mathbf{Y}\right) = E\left(L_0[T^*(\mathbf{Y}) + \epsilon_n] \middle| \mathbf{Y}\right).$$

Hence,

$$R(0, T) \geq E\left[E\left(L_0[T^*(\mathbf{Y}) + \epsilon_n] \middle| \mathbf{Y}\right)\right] = R(0, T^*).$$

□

Corollary 4.1.1 If we take quadratic loss

$$L(\theta, a) := (a - \theta)^2,$$

we get $L_0[a] = a^2$, and so, for $\mathbf{Y} = \mathbf{X} - X_n$,

$$\begin{aligned} T^*(\mathbf{Y}) &= \arg \min_v E\left((v + \epsilon_n)^2 \middle| \mathbf{Y}\right) \\ &= -E(\epsilon_n | \mathbf{Y}), \end{aligned}$$

and hence

$$T^*(\mathbf{X}) = X_n - E(\epsilon_n | \mathbf{Y}).$$

This estimator is called the Pitman estimator.

To investigate the case of quadratic risk further, we:

Note If (X, Z) has density $f(x, z)$ w.r.t. Lebesgue measure, then the density of $Y := X - Z$ is

$$f_Y(y) = \int f(y + z, z) dz.$$

Lemma 4.1.2 Consider quadratic loss. Let \mathbf{p}_0 be the density of $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ w.r.t. Lebesgue measure. Then a UMRE statistic is

$$T^*(\mathbf{X}) = \frac{\int z \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}{\int \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}.$$

Proof. Let $\mathbf{Y} = \mathbf{X} - X_n$. The random vector \mathbf{Y} has density

$$f_{\mathbf{Y}}(y_1, \dots, y_{n-1}, 0) = \int \mathbf{p}_0(y_1 + z, \dots, y_{n-1} + z, z) dz.$$

So the density of ε_n given $\mathbf{Y} = \mathbf{y} = (y_1, \dots, y_{n-1}, 0)$ is

$$f_{\varepsilon_n}(u) = \frac{\mathbf{p}_0(y_1 + u, \dots, y_{n-1} + u, u)}{\int \mathbf{p}_0(y_1 + z, \dots, y_{n-1} + z, z) dz}.$$

It follows that

$$E(\varepsilon_n | \mathbf{y}) = \frac{\int u \mathbf{p}_0(y_1 + u, \dots, y_{n-1} + u, u) du}{\int \mathbf{p}_0(y_1 + z, \dots, y_{n-1} + z, z) dz}.$$

Thus

$$\begin{aligned} E(\varepsilon_n | \mathbf{Y}) &= \frac{\int u \mathbf{p}_0(Y_1 + u, \dots, Y_{n-1} + u, u) du}{\int \mathbf{p}_0(Y_1 + z, \dots, Y_{n-1} + z, z) dz} \\ &= \frac{\int u \mathbf{p}_0(X_1 - X_n + u, \dots, X_{n-1} - X_n + u, u) du}{\int \mathbf{p}_0(X_1 - X_n + z, \dots, X_{n-1} - X_n + z, z) dz} \\ &= \frac{\int u \mathbf{p}_0(X_1 - X_n + u, \dots, X_{n-1} - X_n + u, u) du}{\int \mathbf{p}_0(X_1 + z, \dots, X_{n-1} + z, X_n + z) dz} \\ &= X_n - \frac{\int z \mathbf{p}_0(X_1 - z, \dots, X_{n-1} - z, X_n - z) dz}{\int \mathbf{p}_0(X_1 + z, \dots, X_{n-1} + z, X_n + z) dz}. \end{aligned}$$

Finally, recall that $T^*(\mathbf{X}) = X_n - E(\varepsilon_n | \mathbf{Y})$. □

Example 4.1.1 Suppose X_1, \dots, X_n are i.i.d. Uniform $[\theta - 1/2, \theta + 1/2]$, $\theta \in \mathbb{R}$. Then

$$p_0(x) = \mathbf{1}\{|x| \leq 1/2\}.$$

We have

$$\max_{1 \leq i \leq n} |x_i - z| \leq 1/2 \Leftrightarrow x_{(n)} - 1/2 \leq z \leq x_{(1)} + 1/2.$$

So

$$\mathbf{p}_0(x_1 - z, \dots, x_n - z) = \mathbf{1}\{x_{(n)} - 1/2 \leq z \leq x_{(1)} + 1/2\}.$$

Thus, writing

$$T_1 := X_{(n)} - 1/2, \quad T_2 := X_{(1)} + 1/2,$$

the UMRE estimator T^* is

$$T^* = \left(\int_{T_1}^{T_2} z dz \right) / \left(\int_{T_1}^{T_2} dz \right) = \frac{T_1 + T_2}{2} = \frac{X_{(1)} + X_{(n)}}{2}.$$

We now consider more general invariant statistics \mathbf{Y} .

Definition A map $\mathbf{Y} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called maximal invariant if

$$\mathbf{Y}(\mathbf{x}) = \mathbf{Y}(\mathbf{x}') \Leftrightarrow \exists c : \mathbf{x} = \mathbf{x}' + c.$$

(The constant c may depend on \mathbf{x} and \mathbf{x}' .)

Example The map $\mathbf{Y}(\mathbf{x}) := \mathbf{x} - x_n$ is maximal invariant:

(\Leftarrow) is clear

(\Rightarrow) if $\mathbf{x} - x_n = \mathbf{x}' - x'_n$, we have $\mathbf{x} = \mathbf{x}' + (x_n - x'_n)$.

More generally:

Example Let $d(\mathbf{X})$ be equivariant. Then $\mathbf{Y} := \mathbf{X} - d(\mathbf{X})$ is maximal invariant.

Theorem 4.1.2 Suppose that $d(\mathbf{X})$ is equivariant. Let $\mathbf{Y} := \mathbf{X} - d(\mathbf{X})$, and

$$T^*(\mathbf{Y}) := \arg \min_v E \left(L_0[v + d(\varepsilon)] \middle| \mathbf{Y} \right).$$

Then

$$T^*(\mathbf{X}) := T^*(\mathbf{Y}) + d(\mathbf{X})$$

is UMRE.

Proof. Let T be an equivariant estimator. Then

$$\begin{aligned} T(\mathbf{X}) &= T(\mathbf{X} - d(\mathbf{X})) + d(\mathbf{X}) \\ &= T(\mathbf{Y}) + d(\mathbf{X}). \end{aligned}$$

Hence

$$\begin{aligned} E \left(L_0[T(\varepsilon)] \middle| \mathbf{Y} \right) &= E \left(L_0[T(\mathbf{Y}) + d(\varepsilon)] \middle| \mathbf{Y} \right) \\ &\geq \min_v E \left(L_0[v + d(\varepsilon)] \middle| \mathbf{Y} \right). \end{aligned}$$

Now, use the iterated expectation lemma. □

Special case

For quadratic loss ($L_0[a] = a^2$), the definition of $T^*(\mathbf{Y})$ in the above theorem is

$$T^*(\mathbf{Y}) = -E(d(\varepsilon)|\mathbf{Y}) = -E_0(d(\mathbf{X})|\mathbf{X} - d(\mathbf{X})),$$

so that

$$T^*(\mathbf{X}) = d(\mathbf{X}) - E_0(d(\mathbf{X})|\mathbf{X} - d(\mathbf{X})).$$

So for a equivariant estimator T , we have

$$T \text{ is UMRE} \Leftrightarrow E_0(T(\mathbf{X})|\mathbf{X} - T(\mathbf{X})) = 0.$$

From the right hand side, we conclude that $E_0 T = 0$ and hence $E_\theta(T) = \theta \forall \theta$. Thus, in the case of quadratic loss, an UMRE estimator is unbiased.

Conversely, suppose we have an equivariant and unbiased estimator T . If $T(\mathbf{X})$ and $\mathbf{X} - T(\mathbf{X})$ are independent, it follows that

$$E_0(T(\mathbf{X})|\mathbf{X} - T(\mathbf{X})) = E_0T(\mathbf{X}) = 0.$$

So then T is UMRE.

To check independence, Basu's lemma can be useful.

Basu's lemma *Let X have distribution P_θ , $\theta \in \Theta$. Suppose T is sufficient and complete, and that $Y = Y(X)$ has a distribution that does not depend on θ . Then, for all θ , T and Y are independent under P_θ .*

Proof. Let A be some measurable set, and

$$h(T) := P(Y \in A|T) - P(Y \in A).$$

Notice that indeed, $P(Y \in A|T)$ does not depend on θ because T is sufficient. Because

$$E_\theta h(T) = 0, \quad \forall \theta,$$

we conclude from the completeness of T that

$$h(T) = 0, \quad P_\theta\text{-a.s.}, \quad \forall \theta,$$

in other words,

$$P(Y \in A|T) = P(Y \in A), \quad P_\theta\text{-a.s.}, \quad \forall \theta.$$

Since A was arbitrary, we thus have that the conditional distribution of Y given T is equal to the unconditional distribution:

$$P(Y \in \cdot|T) = P(Y \in \cdot), \quad P_\theta\text{-a.s.}, \quad \forall \theta,$$

that is, for all θ , T and Y are independent under P_θ . □

Basu's lemma is intriguing: it proves a probabilistic property (independence) via statistical concepts.

Example 4.1.2 Let X_1, \dots, X_n be independent $\mathcal{N}(\theta, \sigma^2)$, with σ^2 known. Then $T := \bar{X}$ is sufficient and complete, and moreover, the distribution of $\mathbf{Y} := \mathbf{X} - \bar{X}$ does not depend on θ . So by Basu's lemma, \bar{X} and $\mathbf{X} - \bar{X}$ are independent. Hence, \bar{X} is UMRE.

Remark Indeed, Basu's lemma is peculiar: \bar{X} and $\mathbf{X} - \bar{X}$ of course remain independent if the mean θ is known and/or the variance σ^2 is unknown!

Remark As a by-product, one concludes the independence of \bar{X} and the sample variance $S^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n - 1)$, because S^2 is a function of $\mathbf{X} - \bar{X}$.

4.2 Equivariance in the location-scale model (to be written)

Chapter 5

Proving admissibility and minimaxity

Bayes estimators are quite useful, also for obdurate frequentists. They can be used to construct estimators that are minimax (admissible), or for verification of minimaxity (admissibility).

Let us first recall the definitions. Let $X \in \mathcal{X}$ have distribution P_θ , $\theta \in \Theta$. Let $T = T(X)$ be a statistic (estimator, decision), $L(\theta, a)$ be a loss function, and $R(\theta, T) := E_\theta L(\theta, T(X))$ be the risk of T .

- T is *minimax* if $\forall T' \sup_\theta R(\theta, T) \leq \sup_\theta R(\theta, T')$.
- T is *inadmissible* if $\exists T': \{\forall \theta R(\theta, T') \leq R(\theta, T) \text{ and } \exists \theta R(\theta, T') < R(\theta, T)\}$.
- T is *Bayes* (for the prior density w on Θ) if $\forall T', r_w(T) \leq r_w(T')$.

Recall also that Bayes risk for w is

$$r_w(T) = \int R(\vartheta, T)w(\vartheta)d\mu(\vartheta).$$

Whenever we say that a statistic T is Bayes, without referring to an explicit prior on Θ , we mean that there exists a prior for which T is Bayes. Of course, if the risk $R(\theta, T) = R(T)$ does not depend on θ , then Bayes risk of T does not depend on the prior.

Especially in cases where one wants to use the uniform distribution as prior, but cannot do so because Θ is not bounded, the notion *extended* Bayes is useful.

Definition A statistic T is called *extended* Bayes if there exists a sequence of prior densities $\{w_m\}_{m=1}^\infty$ (*w.r.t. dominating measures that are allowed to depend on m*), such that $r_{w_m}(T) - \inf_{T'} r_{w_m}(T') \rightarrow 0$ as $m \rightarrow \infty$.

5.1 Minimaxity

Lemma 5.1.1 *Suppose T is a statistic with risk $R(\theta, T) = R(T)$ not depending on θ . Then*

(i) T admissible $\Rightarrow T$ minimax,

(ii) T Bayes $\Rightarrow T$ minimax,

and in fact more generally,

(iii) T extended Bayes $\Rightarrow T$ minimax.

Proof.

(i) T is admissible, so for all T' , either there is a θ with $R(\theta, T') > R(T)$, or $R(\theta, T') \geq R(T)$ for all θ . Hence $\sup_{\theta} R(\theta, T') \geq R(T)$.

(ii) Since Bayes implies extended Bayes, this follows from (iii). We nevertheless present a separate proof, as it is somewhat simpler than (iii).

Note first that for any T' ,

$$\begin{aligned} r_w(T') &= \int R(\vartheta, T')w(\vartheta)d\mu(\theta) \leq \int \sup_{\vartheta} R(\vartheta, T')w(\vartheta)d\mu(\theta) \quad (5.1) \\ &= \sup_{\vartheta} R(\vartheta, T'), \end{aligned}$$

that is, Bayes risk is always bounded by the supremum risk. Suppose now that T' is a statistic with $\sup_{\theta} R(\theta, T') < R(T)$. Then

$$r_w(T') \leq \sup_{\vartheta} R(\vartheta, T') < R(T) = r_w(T),$$

which is in contradiction with the assumption that T is Bayes.

(iii) Suppose for simplicity that a Bayes decision T_m for the prior w_m exists, for all m , i.e.

$$r_{w_m}(T_m) = \inf_{T'} r_{w_m}(T'), \quad m = 1, 2, \dots$$

By assumption, for all $\epsilon > 0$, there exists an m sufficiently large, such that

$$R(T) = r_{w_m}(T) \leq r_{w_m}(T_m) + \epsilon \leq r_{w_m}(T') + \epsilon \leq \sup_{\theta} R(\theta, T') + \epsilon,$$

because, as we have seen in (5.1), the Bayes risk is bounded by supremum risk. Since ϵ can be chosen arbitrary small, this proves (iii). \square

Example 5.1.1 Consider a Binomial(n, θ) random variable X . Let the prior on $\theta \in (0, 1)$ be the Beta(r, s) distribution. Then Bayes estimator for quadratic loss is

$$T = \frac{X + r}{n + r + s}.$$

Its risk is

$$\begin{aligned} R(\theta, T) &= E_{\theta}(T - \theta)^2 \\ &= \text{var}_{\theta}(T) + \text{bias}_{\theta}^2(T) \\ &= \frac{n\theta(1 - \theta)}{(n + r + s)^2} + \left[\frac{n\theta + r}{n + r + s} - \frac{(n + r + s)\theta}{n + r + s} \right]^2 \end{aligned}$$

$$= \frac{[(r+s)^2 - n]\theta^2 + [n - 2r(r+s)]\theta + r^2}{(n+r+s)^2}.$$

This can only be constant in θ if the coefficients in front of θ^2 and θ are zero:

$$(r+s)^2 - n = 0, \quad n - 2r(r+s) = 0.$$

Solving for r and s gives

$$r = s = \sqrt{n}/2.$$

Plugging these values back in the estimator T gives

$$T = \frac{X + \sqrt{n}/2}{n + \sqrt{n}}$$

is minimax. The minimax risk is

$$R(T) = \frac{1}{4(\sqrt{n} + 1)^2}.$$

We can compare this with the supremum risk of the unbiased estimator X/n :

$$\sup_{\theta} R(\theta, X/n) = \sup_{\theta} \frac{\theta(1-\theta)}{n} = \frac{1}{4n}.$$

So for large n , this does not differ much from the minimax risk.

Example 5.1.2 We consider again the Pitman estimator (see Lemma 4.1.2)

$$T^* = \int \frac{z \mathbf{p}_0(X_1 - z, \dots, X_n - z)}{\mathbf{p}_0(X_1 - z, \dots, X_n - z)} dz.$$

(Rest is to be written.)

5.2 Admissibility

In this section, the parameter space is assumed to be an open subset of a topological space, so that we can consider open neighborhoods of members of Θ , and continuous functions on Θ . We moreover restrict ourselves to statistics T with $R(\theta, T) < \infty$.

Lemma 5.2.1 *Suppose that the statistic T is Bayes for the prior density w . Then (i) or (ii) below are sufficient conditions for the admissibility of T .*

(i) *The statistic T is the unique Bayes decision (i.e., $r_w(T) = r_w(T')$ implies that $\forall \theta, T = T'$),*

(ii) *For all T' , $R(\theta, T')$ is continuous in θ , and moreover, for all open $U \subset \Theta$, the prior probability $\Pi(U) := \int_U w(\vartheta) d\mu(\vartheta)$ of U is strictly positive.*

Proof.

(i) Suppose that for some T' , $R(\theta, T') \leq R(\theta, T)$ for all θ . Then also $r_w(T') \leq r_w(T)$. Because T is Bayes, we then must have equality:

$$r_w(T') = r_w(T).$$

So then, $\forall \theta$, T' and T are equal P_θ -a.s., and hence, $\forall \theta$, $R(\theta, T') = R(\theta, T)$, so that T' can not be strictly better than T .

(ii) Suppose that T is inadmissible. Then, for some T' , $R(\theta, T') \leq R(\theta, T)$ for all θ , and, for some θ_0 , $R(\theta_0, T') < R(\theta_0, T)$. This implies that for some $\epsilon > 0$, and some open neighborhood $U \subset \Theta$ of θ_0 , we have

$$R(\vartheta, T') \leq R(\vartheta, T) - \epsilon, \quad \vartheta \in U.$$

But then

$$\begin{aligned} r_w(T') &= \int_U R(\vartheta, T')w(\vartheta)d\nu(\vartheta) + \int_{U^c} R(\vartheta, T')w(\vartheta)d\nu(\vartheta) \\ &\leq \int_U R(\vartheta, T)w(\vartheta)d\nu(\vartheta) - \epsilon\Pi(U) + \int_{U^c} R(\vartheta, T)w(\vartheta)d\nu(\vartheta) \\ &= r_w(T) - \epsilon\Pi(U) < r_w(T). \end{aligned}$$

We thus arrived at a contradiction. \square

Lemma 5.2.2 *Suppose that T is extended Bayes, and that for all T' , $R(\theta, T')$ is continuous in θ . In fact assume, for all open sets $U \subset \Theta$,*

$$\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \rightarrow 0,$$

as $m \rightarrow \infty$. Here $\Pi_m(U) := \int_U w_m(\vartheta)d\mu_m(\vartheta)$ is the probability of U under the prior Π_m . Then T is admissible.

Proof. We start out as in the proof of (ii) in the previous lemma. Suppose that T is inadmissible. Then, for some T' , $R(\theta, T') \leq R(\theta, T)$ for all θ , and, for some θ_0 , $R(\theta_0, T') < R(\theta_0, T)$, so that for some $\epsilon > 0$, and some open neighborhood $U \subset \Theta$ of θ_0 , we have

$$R(\vartheta, T') \leq R(\vartheta, T) - \epsilon, \quad \vartheta \in U.$$

This would give that for all m ,

$$r_{w_m}(T') \leq r_{w_m}(T) - \epsilon\Pi_m(U).$$

Suppose for simplicity that a Bayes decision T_m for the prior w_m exists, for all m , i.e.

$$r_{w_m}(T_m) = \inf_{T'} r_{w_m}(T'), \quad m = 1, 2, \dots$$

Then, for all m ,

$$r_{w_m}(T_m) \leq r_{w_m}(T') \leq r_{w_m}(T) - \epsilon\Pi_m(U),$$

or

$$\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \geq \epsilon > 0,$$

that is, we arrived at a contradiction. \square

Example 5.2.1 Let X be $\mathcal{N}(\theta, 1)$ -distributed, and $R(\theta, T) := E_\theta(T - \theta)^2$ be the quadratic risk. We consider estimators of the form

$$T = aX + b, \quad a > 0, \quad b \in \mathbb{R}.$$

Lemma T is admissible if and only if one of the following cases hold

- (i) $a < 1$,
- (ii) $a = 1$ and $b = 0$.

Proof.

(\Leftarrow) (i)

First, we show that T is Bayes for some prior. It turns out that this works with a normal prior, i.e., we take $\theta \sim \mathcal{N}(c, \tau^2)$ for some c and τ^2 to be specified. With the notation

$$f(\vartheta) \propto g(x, \vartheta)$$

we mean that $f(\vartheta)/g(x, \vartheta)$ does not depend on ϑ . We have

$$\begin{aligned} w(\vartheta|x) &= \frac{p(x|\vartheta)w(\vartheta)}{p(x)} \propto \phi(x - \vartheta)\phi\left(\frac{\vartheta - c}{\tau}\right) \\ &\propto \exp\left[-\frac{1}{2}\left\{(x - \vartheta)^2 + \frac{(\vartheta - c)^2}{\tau^2}\right\}\right] \\ &\propto \exp\left[-\frac{1}{2}\left\{\vartheta - \frac{\tau^2 x + c}{\tau^2 + 1}\right\}^2 \frac{1 + \tau^2}{\tau^2}\right]. \end{aligned}$$

We conclude that Bayes estimator is

$$T_{\text{Bayes}} = E(\theta|X) = \frac{\tau^2 X + c}{\tau^2 + 1}.$$

Taking

$$\frac{\tau^2}{\tau^2 + 1} = a, \quad \frac{c}{\tau^2 + 1} = b,$$

yields $T = T_{\text{Bayes}}$.

Next, we check (i) in Lemma 5.2.1, i.e. that T is unique. For quadratic loss, and for $T = E(\theta|X)$, the Bayes risk of an estimator T' is

$$r_w(T') = E\text{var}(\theta|X) + E(T - T')^2.$$

This follows from straightforward calculations:

$$\begin{aligned} r_w(T') &= \int R(\vartheta, T')w(\vartheta)d\mu(\vartheta) \\ &= ER(\theta, T') = E(\theta - T')^2 = E\left[E\left((\theta - T')^2 \middle| X\right)\right] \end{aligned}$$

and, with θ being the random variable,

$$E\left((\theta - T')^2 \middle| X\right) = E\left((\theta - T)^2 \middle| X\right) + (T - T')^2 = \text{var}(\theta|X) + (T - T')^2.$$

We conclude that if $r_w(T') = r_w(T)$, then

$$E(T - T')^2 = 0.$$

Here, the expectation is with θ integrated out, i.e., with respect to the measure P with density

$$p(x) = \int p_\vartheta(x)w(\vartheta)d\mu(\vartheta).$$

Now, we can write $X = \theta + \epsilon$, with $\theta \mathcal{N}(c, \tau^2)$ -distributed, and with ϵ a standard normal random variable independent of θ . So X is $\mathcal{N}(c, \tau^2 + 1)$, that is, P is the $\mathcal{N}(c, \tau^2 + 1)$ -distribution. Now, $E(T - T')^2 = 0$ implies $T = T'$ P -a.s.. Since P dominates all P_θ , we conclude that $T = T'$ P_θ -a.s., for all θ . So T is unique, and hence admissible.

(\Leftarrow) (ii)

In this case, $T = X$. We use Lemma 5.2.2. Because $R(\theta, T) = 1$ for all θ , also $r_w(T) = 1$ for any prior. Let w_m be the density of the $\mathcal{N}(0, m)$ -distribution. As we have seen in the previous part of the proof, the Bayes estimator is

$$T_m = \frac{m}{m+1}X.$$

By the bias-variance decomposition, it has risk

$$R(\theta, T_m) = \frac{m^2}{(m+1)^2} + \left(\frac{m}{m+1} - 1\right)^2 \theta^2 = \frac{m^2}{(m+1)^2} + \frac{\theta^2}{(m+1)^2}.$$

As $E\theta^2 = m$, its Bayes risk is

$$r_{w_m}(T_m) = \frac{m^2}{(m+1)^2} + \frac{m}{(m+1)^2} = \frac{m}{m+1}.$$

It follows that

$$r_{w_m}(T) - r_{w_m}(T_m) = 1 - \frac{m}{m+1} = \frac{1}{m+1}.$$

So T is extended Bayes. But we need to prove the more refined property of Lemma 5.2.2. It is clear that here, we only need to consider open intervals $U = (u, u+h)$, with u and $h > 0$ fixed. We have

$$\begin{aligned} \Pi_m(U) &= \Phi\left(\frac{u+h}{\sqrt{m}}\right) - \Phi\left(\frac{u}{\sqrt{m}}\right) \\ &= \frac{1}{\sqrt{m}}\phi\left(\frac{u}{\sqrt{m}}\right)h + o(1/\sqrt{m}). \end{aligned}$$

For m large,

$$\phi\left(\frac{u}{\sqrt{m}}\right) \approx \phi(0) = \frac{1}{\sqrt{2\pi}} > \frac{1}{4} \text{ (say),}$$

so for m sufficiently large (depending on u)

$$\phi\left(\frac{u}{\sqrt{m}}\right) \geq \frac{1}{4}.$$

Thus, for m sufficiently large (depending on u and h), we have

$$\Pi_m(U) \geq \frac{1}{4\sqrt{m}}h.$$

We conclude that for m sufficiently large

$$\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \leq \frac{4}{h\sqrt{m}}.$$

As the right hand side converges to zero as $m \rightarrow \infty$, this shows that X is admissible.

(\Rightarrow)

We now have to show that if (i) or (ii) do not hold, then T is not admissible. This means we have to consider two cases: $a > 1$ and $a = 1, b \neq 0$. In the case $a > 1$, we have $R(\theta, aX + b) \geq \text{var}(aX + b) > 1 = R(\theta, X)$, so $aX + b$ is not admissible. When $a = 1$ and $b \neq 0$, it is the bias term that makes $aX + b$ inadmissible:

$$R(\theta, X + b) = 1 + b^2 > 1 = R(\theta, X).$$

.

□

Lemma 5.2.3 *Let $\theta \in \Theta = \mathbb{R}$ and $\{P_\theta : \theta \in \Theta\}$ be an exponential family in canonical form:*

$$p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x).$$

Then T is an admissible estimator of $g(\theta) := \dot{d}(\theta)$, under quadratic loss (i.e., under the loss $L(\theta, a) := |a - g(\theta)|^2$).

Proof. Recall that

$$\dot{d}(\theta) = E_\theta T, \quad \ddot{d}(\theta) = \text{var}_\theta(T) = I(\theta).$$

Now, let T' be some estimator, with expectation

$$E_\theta T' := q(\theta).$$

the bias of T' is

$$b(\theta) = q(\theta) - g(\theta),$$

or

$$q(\theta) = b(\theta) + g(\theta) = b(\theta) + \dot{d}(\theta).$$

This implies

$$\dot{q}(\theta) = \dot{b}(\theta) + I(\theta).$$

By the Cramer Rao lower bound

$$\begin{aligned} R(\theta, T') &= \text{var}_\theta(T') + b^2(\theta) \\ &\geq \frac{[\dot{q}(\theta)]^2}{I(\theta)} + b^2(\theta) = \frac{[\dot{b}(\theta) + I(\theta)]^2}{I(\theta)} + b^2(\theta). \end{aligned}$$

Suppose now that

$$R(\theta, T') \leq R(\theta, T), \forall \theta.$$

Because $R(\theta, T) = I(\theta)$ this implies

$$\frac{[\dot{b}(\theta) + I(\theta)]^2}{I(\theta)} + b^2(\theta) \leq I(\theta),$$

or

$$I(\theta)\{b^2(\theta) + 2\dot{b}(\theta)\} \leq -[\dot{b}(\theta)]^2 \leq 0.$$

This in turn implies

$$b^2(\theta) + 2\dot{b}(\theta) \leq 0,$$

end hence

$$\frac{\dot{b}(\theta)}{b^2(\theta)} \leq -\frac{1}{2},$$

so

$$\frac{d}{d\theta} \left(\frac{1}{b(\theta)} \right) - \frac{1}{2} \geq 0,$$

or

$$\frac{d}{d\theta} \left(\frac{1}{b(\theta)} - \frac{\theta}{2} \right) \geq 0.$$

In other words, $1/b(\theta) - \theta/2$ is an increasing function, and hence, $b(\theta)$ is a decreasing function.

We will now show that this implies that $b(\theta) = 0$ for all θ .

Suppose instead $b(\theta_0) < 0$ for some θ_0 . Let $\theta_0 < \vartheta \rightarrow \infty$. Then

$$\frac{1}{b(\vartheta)} \geq \frac{1}{b(\theta_0)} + \frac{\vartheta - \theta_0}{2} \rightarrow \infty,$$

i.e.,

$$b(\vartheta) \rightarrow 0, \quad \vartheta \rightarrow \infty$$

This is not possible, as $b(\theta)$ is a decreasing function.

Similarly, if $b(\theta_0) > 0$, take $\theta_0 \geq \vartheta \rightarrow -\infty$, to find again

$$b(\vartheta) \rightarrow 0, \quad \vartheta \rightarrow -\infty,$$

which is not possible.

We conclude that $b(\theta) = 0$ for all θ , i.e., T' is an unbiased estimator of θ . By the Cramer Rao lower bound, we now conclude

$$R(\theta, T') = \text{var}_\theta(T') \geq R(\theta, T) = I(\theta).$$

□

Example Let X be $\mathcal{N}(\theta, 1)$ -distributed, with $\theta \in \mathbb{R}$ unknown. Then X is an admissible estimator of θ .

Example Let X be $\mathcal{N}(0, \sigma^2)$, with $\sigma^2 \in (0, \infty)$ unknown. Its density is

$$p_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{x^2}{2\sigma^2}\right] = \exp[\theta T(x) - d(\theta)]h(x),$$

with

$$T(x) = -x^2/2, \quad \theta = 1/\sigma^2, \quad d(\theta) = (\log \sigma^2)/2 = -(\log \theta)/2,$$

$$\dot{d}(\theta) = -\frac{1}{2\theta} = -\frac{\sigma^2}{2},$$

$$\ddot{d}(\theta) = \frac{1}{2\theta^2} = \frac{\sigma^4}{2}.$$

Observe that $\theta \in \Theta = (0, \infty)$, which is not the whole real line. So Lemma 5.2.3 cannot be applied. We will now show that T is not admissible. Define for all $a > 0$,

$$T_a := -aX^2.$$

so that $T = T_{1/2}$. We have

$$R(\theta, T_a) = \text{var}_\theta(T_a) + \text{bias}_\theta^2(T_a)$$

$$= 2a^2\sigma^4 + [a - 1/2]^2\sigma^4.$$

Thus, $R(\theta, T_a)$ is minimized at $a = 1/6$ giving

$$R(\theta, T_{1/6}) = \sigma^4/6 < \sigma^4/2 = R(\theta, T).$$

5.3 Inadmissibility in higher-dimensional settings (to be written)

Chapter 6

Asymptotic theory

In this chapter, the observations X_1, \dots, X_n are considered as the first n of an infinite sequence of i.i.d. random variables X_1, \dots, X_n, \dots with values in \mathcal{X} and with distribution P . We say that the X_i are i.i.d. *copies*, of some random variable $X \in \mathcal{X}$ with distribution P . We let $\mathbb{P} = P \times P \times \dots$ be the distribution of the whole sequence $\{X_i\}_{i=1}^\infty$.

The model class for P is $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$. When $P = P_\theta$, we write $\mathbb{P} = \mathbb{P}_\theta = P_\theta \times P_\theta \times \dots$. The parameter of interest is

$$\gamma := g(\theta) \in \mathbb{R}^p,$$

where $g : \Theta \rightarrow \mathbb{R}^p$ is a given function. We let

$$\Gamma := \{g(\theta) : \theta \in \Theta\}$$

be the parameter space for γ .

An estimator of γ , based on the data X_1, \dots, X_n , is some function $T_n = T_n(X_1, \dots, X_n)$ of the data. We assume the estimator is defined for all n , i.e., we actually consider a sequence of estimators $\{T_n\}_{n=1}^\infty$.

Remark Under the i.i.d. assumption, it is natural to assume that each T_n is a symmetric function of the data, that is

$$T_n(X_1, \dots, X_n) = T_n(X_{\pi_1}, \dots, X_{\pi_n})$$

for all permutations π of $\{1, \dots, n\}$. In that case, one can write T_n in the form $T_n = Q(\hat{P}_n)$, where \hat{P}_n is the empirical distribution (see also Subsection 1.9.1).

6.1 Types of convergence

Definition Let $\{Z_n\}_{n=1}^\infty$ and Z be \mathbb{R}^p -valued random variables defined on the same probability space.¹ We say that Z_n converges in probability to Z if for all

¹Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space, and $X : \Omega \rightarrow \mathcal{X}$ and $Y : \Omega \rightarrow \mathcal{Y}$ be two measurable maps. Then X and Y are called random variables, and they are defined on the same probability space Ω .

$\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\|Z_n - Z\| > \epsilon) = 0.$$

Notation: $Z_n \xrightarrow{\mathbb{P}} Z$.

Remark Chebyshev's inequality can be a tool to prove convergence in probability. It says that for all increasing functions $\psi : [0, \infty) \rightarrow [0, \infty)$, one has

$$\mathbb{P}(\|Z_n - Z\| \geq \epsilon) \leq \frac{\mathbb{E}\psi(\|Z_n - Z\|)}{\psi(\epsilon)}.$$

Definition Let $\{Z_n\}_{n=1}^{\infty}$ and Z be \mathbb{R}^p -valued random variables. We say that Z_n converges in distribution to Z , if for all continuous and bounded functions f ,

$$\lim_{n \rightarrow \infty} \mathbb{E}f(Z_n) = \mathbb{E}f(Z).$$

Notation: $Z_n \xrightarrow{D} Z$.

Remark Convergence in probability implies convergence in distribution, but not the other way around.

Example Let X_1, X_2, \dots be i.i.d. real-valued random variables with mean μ and variance σ^2 . Let $\bar{X}_n := \sum_{i=1}^n X_i/n$ be the average of the first n . Then by the central limit theorem (CLT),

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \mathcal{N}(0, \sigma^2),$$

that is

$$\mathbb{P}\left(\sqrt{n} \frac{(\bar{X}_n - \mu)}{\sigma} \leq z\right) \rightarrow \Phi(z), \quad \forall z.$$

The following theorem says that for convergence in distribution, one actually can do with one-dimensional random variables. We omit the proof.

Theorem 6.1.1 (*Cramér-Wold device*) Let $(\{Z_n\}, Z)$ be a collection of \mathbb{R}^p -valued random variables. Then

$$Z_n \xrightarrow{D} Z \Leftrightarrow a^T Z_n \xrightarrow{D} a^T Z \quad \forall a \in \mathbb{R}^p.$$

Example Let X_1, X_2, \dots be i.i.d. copies of a random variable $X = (X^{(1)}, \dots, X^{(p)})^T$ in \mathbb{R}^p . Assume $EX := \mu = (\mu_1, \dots, \mu_p)^T$ and $\Sigma := \text{Cov}(X) := EXX^T - \mu\mu^T$ exist. Then for all $a \in \mathbb{R}^p$,

$$Ea^T X = a^T \mu, \quad \text{var}(a^T X) = a^T \Sigma a.$$

Define

$$\bar{X}_n = (\bar{X}_n^{(1)}, \dots, \bar{X}_n^{(p)})^T.$$

By the 1-dimensional CLT, for all $a \in \mathbb{R}^p$,

$$\sqrt{n}(a^T \bar{X}_n - a^T \mu) \xrightarrow{D} \mathcal{N}(0, a^T \Sigma a).$$

The Cramér-Wold device therefore gives the p -dimensional CLT

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \mathcal{N}(0, \Sigma).$$

We recall the *Portmanteau Theorem*:

Theorem 6.1.2 *Let $(\{Z_n\}, Z)$ be a collection of \mathbb{R}^p -valued random variables. Denote the distribution of Z by Q and let $G = Q(Z \leq \cdot)$ be its distribution function. The following statements are equivalent:*

- (i) $Z_n \xrightarrow{D} Z$ (i.e., $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z) \forall f$ bounded and continuous).
- (ii) $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z) \forall f$ bounded and Lipschitz.²
- (iii) $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z) \forall f$ bounded and Q -a.s. continuous.
- (iv) $\mathbb{P}(Z_n \leq z) \rightarrow G(z)$ for all G -continuity points z .

6.1.1 Stochastic order symbols

Let $\{Z_n\}$ be a collection of \mathbb{R}^p -valued random variables, and let $\{r_n\}$ be strictly positive random variables. We write

$$Z_n = O_{\mathbf{P}}(1)$$

(Z_n is bounded in probability) if

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(\|Z_n\| > M) = 0.$$

This is also called *uniform tightness* of the sequence $\{Z_n\}$. We write $Z_n = O_{\mathbf{P}}(r_n)$ if $Z_n/r_n = O_{\mathbf{P}}(1)$.

If Z_n converges in probability to zero, we write this as

$$Z_n = o_{\mathbf{P}}(1).$$

Moreover, $Z_n = o_{\mathbf{P}}(r_n)$ (Z_n is of small order r_n in probability) if $Z_n/r_n = o_{\mathbf{P}}(1)$.

6.1.2 Some implications of convergence

Lemma 6.1.1 *Suppose that Z_n converges in distribution. Then $Z_n = O_{\mathbf{P}}(1)$.*

²A real-valued function f on (a subset of) \mathbb{R}^p is *Lipschitz* if for a constant C and all (z, \tilde{z}) in the domain of f , $|f(z) - f(\tilde{z})| \leq C\|z - \tilde{z}\|$.

Proof. To simplify, take $p = 1$ (Cramér-Wold device). Let $Z_n \xrightarrow{D} Z$, where Z has distribution function G . Then for every G -continuity point M ,

$$\mathbb{P}(Z_n > M) \rightarrow 1 - G(M),$$

and for every G -continuity point $-M$,

$$\mathbb{P}(Z_n \leq -M) \rightarrow G(-M).$$

Since $1 - G(M)$ as well as $G(-M)$ converge to zero as $M \rightarrow \infty$, the result follows. \square

Example Let X_1, X_2, \dots be i.i.d. copies of a random variable $X \in \mathbb{R}$ with $EX = \mu$ and $\text{var}(X) < \infty$. Then by the CLT,

$$\bar{X}_n - \mu = O_{\mathbf{P}}\left(\frac{1}{\sqrt{n}}\right).$$

Theorem 6.1.3 (Slutsky) Let $(\{Z_n, A_n\}, Z)$ be a collection of \mathbb{R}^p -valued random variables, and $a \in \mathbb{R}^p$ be a vector of constants. Assume that $Z_n \xrightarrow{D} Z$, $A_n \xrightarrow{\mathbf{P}} a$. Then

$$A_n^T Z_n \xrightarrow{D} a^T Z.$$

Proof. Take a bounded Lipschitz function f , say

$$|f| \leq C_B, \quad |f(z) - f(\tilde{z})| \leq C_L \|z - \tilde{z}\|.$$

Then

$$\begin{aligned} & \left| \mathbf{E}f(A_n^T Z_n) - \mathbf{E}f(a^T Z) \right| \\ & \leq \left| \mathbf{E}f(A_n^T Z_n) - \mathbf{E}f(a^T Z_n) \right| + \left| \mathbf{E}f(a^T Z_n) - \mathbf{E}f(a^T Z) \right|. \end{aligned}$$

Because the function $z \mapsto f(a^T z)$ is bounded and Lipschitz (with Lipschitz constant $\|a\|C_L$), we know that the second term goes to zero. As for the first term, we argue as follows. Let $\epsilon > 0$ and $M > 0$ be arbitrary. Define $S_n := \{\|Z_n\| \leq M, \|A_n - a\| \leq \epsilon\}$. Then

$$\begin{aligned} & \left| \mathbf{E}f(A_n^T Z_n) - \mathbf{E}f(a^T Z_n) \right| \leq \mathbf{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \\ & = \mathbf{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \mathbf{1}_{\{S_n\}} + \mathbf{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \mathbf{1}_{\{S_n^c\}} \\ & \leq C_L \epsilon M + 2C_B \mathbf{P}(S_n^c). \end{aligned} \tag{6.1}$$

Now

$$\mathbf{P}(S_n^c) \leq \mathbf{P}(\|Z_n\| > M) + \mathbf{P}(\|A_n - a\| > \epsilon).$$

Thus, both terms in (6.1) can be made arbitrary small by appropriately choosing ϵ small and n and M large. \square

6.2 Consistency and asymptotic normality

Definition A sequence of estimators $\{T_n\}$ of $\gamma = g(\theta)$ is called consistent if

$$T_n \xrightarrow{\mathbb{P}_\theta} \gamma.$$

Definition A sequence of estimators $\{T_n\}$ of $\gamma = g(\theta)$ is called asymptotically normal with asymptotic covariance matrix V_θ , if

$$\sqrt{n}(T_n - \gamma) \xrightarrow{D_\theta} \mathcal{N}(0, V_\theta).$$

Example Suppose \mathcal{P} is the location model

$$\mathcal{P} = \{P_{\mu, F_0}(X \leq \cdot) := F_0(\cdot - \mu), \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0\}.$$

The parameter is then $\theta = (\mu, F_0)$ and $\Theta = \mathbb{R} \times \mathcal{F}_0$. We assume for all $F_0 \in \mathcal{F}_0$

$$\int x dF_0(x) = 0, \quad \sigma_{F_0}^2 := \int x^2 dF_0(x) < \infty.$$

Let $g(\theta) := \mu$ and $T_n := (X_1 + \cdots + X_n)/n = \bar{X}_n$. Then T_n is a consistent estimator of μ and, by the central limit theorem

$$\sqrt{n}(T_n - \mu) \xrightarrow{D_\theta} \mathcal{N}(0, \sigma_{F_0}^2).$$

6.2.1 Asymptotic linearity

As we will show, for many estimators, asymptotic normality is a consequence of asymptotic linearity, that is, the estimator is approximately an average, to which we can apply the CLT.

Definition The sequence of estimators $\{T_n\}$ of $\gamma = g(\theta)$ is called asymptotically linear if for a function $l_\theta : \mathcal{X} \rightarrow \mathbb{R}^p$, with $E_\theta l_\theta(X) = 0$ and

$$E_\theta l_\theta(X) l_\theta^T(X) := V_\theta < \infty,$$

it holds that

$$T_n - \gamma = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbb{P}_\theta}(n^{-1/2}).$$

Remark. We then call l_θ the *influence function* of (the sequence) T_n . Roughly speaking, $l_\theta(x)$ approximately measures the influence of an additional observation x .

Example Assuming the entries of X have finite variance, the estimator $T_n := \bar{X}_n$ is a linear and hence asymptotically linear estimator of the mean μ , with influence function

$$l_\theta(x) = x - \mu.$$

Example 6.2.1 Let X be real-valued, with $E_\theta X := \mu$, $\text{var}_\theta(X) := \sigma^2$ and $\kappa := E_\theta(X - \mu)^4$ (assumed to exist). Consider the estimator

$$\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

of σ^2 . We rewrite

$$\begin{aligned} \hat{\sigma}_n^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + (\bar{X}_n - \mu)^2 - \frac{2}{n} \sum_{i=1}^n (X_i - \mu)(\bar{X}_n - \mu) \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X}_n - \mu)^2. \end{aligned}$$

Because by the CLT, $\bar{X}_n - \mu = O_{\mathbf{P}_\theta}(n^{-1/2})$, we get

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + O_{\mathbf{P}_\theta}(1/n).$$

So $\hat{\sigma}_n^2$ is asymptotically linear with influence function

$$l_\theta(x) = (x - \mu)^2 - \sigma^2.$$

The asymptotic variance is

$$V_\theta = E_\theta \left((X - \mu)^2 - \sigma^2 \right)^2 = \kappa - \sigma^4.$$

6.2.2 The δ -technique

Theorem 6.2.1 Let $(\{T_n\}, Z)$ be a collection of random variables in \mathbb{R}^p , $c \in \mathbb{R}^p$ be a nonrandom vector, and $\{r_n\}$ be a nonrandom sequence of positive numbers, with $r_n \downarrow 0$. Moreover, let $h : \mathbb{R}^p \rightarrow \mathbb{R}$ be differentiable at c , with derivative $\dot{h}(c) \in \mathbb{R}^p$. Suppose that

$$(T_n - c)/r_n \xrightarrow{\mathbf{D}} Z.$$

Then

$$(h(T_n) - h(c))/r_n \xrightarrow{\mathbf{D}} \dot{h}(c)^T Z.$$

Proof. By Slutsky's Theorem,

$$\dot{h}(c)^T (T_n - c)/r_n \xrightarrow{\mathbf{D}} \dot{h}(c)^T Z.$$

Since $(T_n - c)/r_n$ converges in distribution, we know that $\|T_n - c\|/r_n = O_{\mathbf{P}}(1)$. Hence, $\|T_n - c\| = O_{\mathbf{P}}(r_n)$. The result follows now from

$$h(T_n) - h(c) = \dot{h}(c)^T (T_n - c) + o(\|T_n - c\|) = \dot{h}(c)^T (T_n - c) + o_{\mathbf{P}}(r_n).$$

□

Corollary 6.2.1 Let T_n be an asymptotically linear estimator of $\gamma := g(\theta)$, with influence function l_θ and asymptotic covariance matrix V_θ . Suppose h is differentiable at γ . Then it follows in the same way as in the previous theorem, that $h(T_n)$ is an asymptotically linear estimator of $h(\gamma)$, with influence function $\dot{h}(\gamma)^T l_\theta$ and asymptotic variance $\dot{h}(\gamma)^T V_\theta \dot{h}(\gamma)$.

Example 6.2.2 Let X_1, \dots, X_n be a sample from the Exponential(θ) distribution, with $\theta > 0$. Then \bar{X}_n is a linear estimator of $E_\theta X = 1/\theta := \gamma$, with influence function $l_\theta(x) = x - 1/\theta$. The variance of $\sqrt{n}(T_n - 1/\theta)$ is $1/\theta^2 = \gamma^2$. Thus, $1/\bar{X}_n$ is an asymptotically linear estimator of θ . In this case, $h(\gamma) = 1/\gamma$, so that $\dot{h}(\gamma) = -1/\gamma^2$. The influence function of $1/\bar{X}_n$ is thus

$$\dot{h}(\gamma)l_\theta(x) = -\frac{1}{\gamma^2}(x - \gamma) = -\theta^2(x - 1/\theta).$$

The asymptotic variance of $1/\bar{X}_n$ is

$$[\dot{h}(\gamma)]^2 \gamma^2 = \frac{1}{\gamma^2} = \theta^2.$$

So

$$\sqrt{n} \left(\frac{1}{\bar{X}_n} - \theta \right) \xrightarrow{D_\theta} \mathcal{N}(0, \theta^2).$$

Example 6.2.3 Consider again Example 6.2.1. Let X be real-valued, with $E_\theta X := \mu$, $\text{var}_\theta(X) := \sigma^2$ and $\kappa := E_\theta(X - \mu)^4$ (assumed to exist). Define moreover, for $r = 1, 2, 3, 4$, the r -th moment $\mu_r := E_\theta X^r$. We again consider the estimator

$$\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

We have

$$\hat{\sigma}_n^2 = h(T_n),$$

where $T_n = (T_{n,1}, T_{n,2})^T$, with

$$T_{n,1} = \bar{X}_n, \quad T_{n,2} = \frac{1}{n} \sum_{i=1}^n X_i^2,$$

and

$$h(t) = t_2 - t_1^2, \quad t = (t_1, t_2)^T.$$

The estimator T_n has influence function

$$l_\theta(x) = \begin{pmatrix} x - \mu_1 \\ x^2 - \mu_2 \end{pmatrix}.$$

By the 2-dimensional CLT,

$$\sqrt{n} \left(T_n - \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \right) \xrightarrow{D_\theta} \mathcal{N}(0, \Sigma),$$

with

$$\Sigma = \begin{pmatrix} \mu_2 - \mu_1^2 & \mu_3 - \mu_1\mu_2 \\ \mu_3 - \mu_1\mu_2 & \mu_4 - \mu_2^2 \end{pmatrix}.$$

It holds that

$$\dot{h}\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}\right) = \begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix},$$

so that $\hat{\sigma}_n^2$ has influence function

$$\begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix}^T \begin{pmatrix} x - \mu_1 \\ x^2 - \mu_2 \end{pmatrix} = (x - \mu)^2 - \sigma^2,$$

(invoking $\mu_1 = \mu$). After some calculations, one finds moreover that

$$\begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix}^T \Sigma \begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix} = \kappa - \sigma^4,$$

i.e., the δ -method gives the same result as the ad hoc method in Example 6.2.1, as it of course should.

6.3 M-estimators

Let, for each $\gamma \in \Gamma$, be defined some loss function $\rho_\gamma(X)$. These are for instance constructed as in Chapter 2: we let $L(\theta, a)$ be the loss when taking action a . Then, we fix some decision $d(x)$, and rewrite

$$L(\theta, d(x)) := \rho_\gamma(x),$$

assuming the loss L depends only on θ via the parameter of interest $\gamma = g(\theta)$.

We now require that the risk

$$E_\theta \rho_c(X)$$

is minimized at the value $c = \gamma$ i.e.,

$$\gamma = \arg \min_{c \in \Gamma} E_\theta \rho_c(X). \quad (6.2)$$

Alternatively, given ρ_c , one may view (6.2) as the *definition* of γ .

If $c \mapsto \rho_c(x)$ is differentiable for all x , we write

$$\psi_c(x) := \dot{\rho}_c(x) := \frac{\partial}{\partial c} \rho_c(x).$$

Then, assuming we may interchange differentiation and taking expectations³, we have

$$E_\theta \psi_\gamma(X) = 0.$$

³If $|\partial \rho_c / \partial c| \leq H(\cdot)$ where $E_\theta H(X) < \infty$, then it follows from the dominated convergence theorem that $\partial[E_\theta \rho_c(X)] / \partial c = E_\theta[\partial \rho_c(X) / \partial c]$.

Example 6.3.1 Let $X \in \mathbb{R}$, and let the parameter of interest be the mean $\mu = E_\theta X$. Assume X has finite variance σ^2 . Then

$$\mu = \arg \min_c E_\theta (X - c)^2,$$

as (recall), by the bias-variance decomposition

$$E_\theta (X - c)^2 = \sigma^2 + (\mu - c)^2.$$

So in this case, we can take

$$\rho_c(x) = (x - c)^2.$$

Example 6.3.2 Suppose $\Theta \subset \mathbb{R}^p$ and that the densities $p_\theta = dP_\theta/d\nu$ exist w.r.t. some σ -finite measure ν .

Definition *The quantity*

$$K(\tilde{\theta}|\theta) = E_\theta \log \left(\frac{p_{\tilde{\theta}}(X)}{p_\theta(X)} \right)$$

is called the Kullback Leibler information, or the relative entropy.

Remark Some care has to be taken, not to divide by zero! This can be handled e.g., by assuming that the support $\{x : p_\theta(x) > 0\}$ does not depend on θ (see also condition I in the CRLB of Chapter 3).



Define now

$$\rho_\theta(x) = -\log p_\theta(x).$$

One easily sees that

$$K(\tilde{\theta}|\theta) = E_\theta \rho_{\tilde{\theta}}(X) - E_\theta \rho_\theta(X).$$

Lemma $E_\theta \rho_{\tilde{\theta}}(X)$ is minimized at $\tilde{\theta} = \theta$:

$$\theta = \arg \min_{\tilde{\theta}} E_\theta \rho_{\tilde{\theta}}(X).$$

Proof. We will show that

$$K(\tilde{\theta}|\theta) \geq 0.$$

This follows from Jensen's inequality. Since the log-function is concave,

$$K(\tilde{\theta}|\theta) = -E_\theta \log \left(\frac{p_{\tilde{\theta}}(X)}{p_\theta(X)} \right) \geq -\log \left(E_\theta \left(\frac{p_{\tilde{\theta}}(X)}{p_\theta(X)} \right) \right) = -\log 1 = 0.$$

□

Definition The M-estimator $\hat{\gamma}_n$ of γ is defined as

$$\hat{\gamma}_n := \arg \min_{c \in \Gamma} \frac{1}{n} \sum_{i=1}^n \rho_c(X_i).$$

The “M” in “M-estimator” stands for Minimizer (or - take minus signs - Maximizer).

If $\rho_c(x)$ is differentiable in c for all x , we generally can define $\hat{\gamma}_n$ as the solution of putting the derivatives

$$\frac{\partial}{\partial c} \sum_{i=1}^n \rho_c(X_i) = \sum_{i=1}^n \psi_c(X_i)$$

to zero. This is called the Z-estimator.

Definition The Z-estimator $\hat{\gamma}_n$ of γ is defined as a solution of the equations

$$\frac{1}{n} \sum_{i=1}^n \psi_{\hat{\gamma}_n}(X_i) = 0.$$

Remark A solution $\hat{\gamma}_n \in \Gamma$ is then assumed to exist.

6.3.1 Consistency of M-estimators

Note that γ minimizes a theoretical expectation, whereas the M-estimator $\hat{\gamma}_n$ minimizes the empirical average. Likewise, γ is a solution of putting a theoretical expectation to zero, whereas the Z-estimator $\hat{\gamma}_n$ is the solution of putting an empirical average to zero.

By the law of large numbers, averages converge to expectations. So the M-estimator (Z-estimator) does make sense. However, consistency and further properties are not immediate, because we actually need convergence the averages to expectations over a range of values $c \in \Gamma$ simultaneously. This is the topic of *empirical process theory*.

We will borrow the notation from empirical process theory. That is, for a function $f : \mathcal{X} \rightarrow \mathbb{R}^r$, we let

$$P_\theta f := E_\theta f(X), \quad \hat{P}_n f := \frac{1}{n} \sum_{i=1}^n f(X_i).$$

Then, by the law of large numbers, if $P_\theta |f| < \infty$,

$$(\hat{P}_n - P_\theta)f \rightarrow 0, \quad \mathbf{P}_\theta\text{-a.s.}$$

We now show consistency assuming the parameter space Γ is compact. (The assumption of compactness can however often be omitted if $c \mapsto \rho_c$ is convex. We skip the details.)

We will need that convergence of to the minimum value also implies convergence of the arg min, i.e., convergence of the location of the minimum. To this end, we assume

Definition *The minimizer γ of $P_\theta \rho_c$ is called well-separated if for all $\epsilon > 0$,*

$$\inf\{P_\theta \rho_c : c \in \Gamma, \|c - \gamma\| > \epsilon\} > P_\theta \rho_\gamma.$$

Theorem 6.3.1 *Suppose that Γ is compact, that $c \mapsto \rho_c(x)$ is continuous for all x , and that*

$$P_\theta \left(\sup_{c \in \Gamma} |\rho_c| \right) < \infty.$$

Then

$$P_\theta \rho_{\hat{\gamma}_n} \rightarrow P_\theta \rho_\gamma, \mathbf{P}_\theta\text{-a.s..}$$

If γ is well-separated, this implies $\hat{\gamma}_n \rightarrow \gamma$, $\mathbf{P}_\theta\text{-a.s..}$

Proof. We will show the uniform convergence

$$\sup_{c \in \Gamma} |(\hat{P}_n - P_\theta) \rho_c| \rightarrow 0, P_\theta\text{-a.s..} \quad (6.3)$$

This will indeed suffice, as

$$\begin{aligned} 0 &\leq P_\theta(\rho_{\hat{\gamma}_n} - \rho_\gamma) = -(\hat{P}_n - P_\theta)(\rho_{\hat{\gamma}_n} - \rho_\gamma) + \hat{P}_n(\rho_{\hat{\gamma}_n} - \rho_\gamma) \\ &\leq -(\hat{P}_n - P_\theta)(\rho_{\hat{\gamma}_n} - \rho_\gamma) \leq |(\hat{P}_n - P_\theta)\rho_{\hat{\gamma}_n}| + |(\hat{P}_n - P_\theta)\rho_\gamma| \\ &\leq \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta)\rho_c| + |(\hat{P}_n - P_\theta)\rho_\gamma| \leq 2 \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta)\rho_c|. \end{aligned}$$

To show (6.3), define for each $\delta > 0$ and $c \in \Gamma$,

$$w(\cdot, \delta, c) := \sup_{\tilde{c} \in \Gamma: \|\tilde{c} - c\| < \delta} |\rho_{\tilde{c}} - \rho_c|.$$

Then for all x , as $\delta \downarrow 0$,

$$w(x, \delta, c) \rightarrow 0.$$

So also, by dominated convergence

$$P_\theta w(\cdot, \delta, c) \rightarrow 0.$$

Hence, for all $\epsilon > 0$, there exists a δ_c such that

$$P_\theta w(\cdot, \delta_c, c) \leq \epsilon.$$

Let

$$B_c := \{\tilde{c} \in \Gamma : \|\tilde{c} - c\| < \delta_c\}.$$

Then $\{B_c : c \in \Gamma\}$ is a covering of Γ by open sets. Since Γ is compact, there exists finite sub-covering

$$B_{c_1} \dots B_{c_N}.$$

For $c \in B_{c_j}$,

$$|\rho_c - \rho_{c_j}| \leq w(\cdot, \delta_{c_j}, c_j).$$

It follows that

$$\begin{aligned} \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta)\rho_c| &\leq \max_{1 \leq j \leq N} |(\hat{P}_n - P_\theta)\rho_{c_j}| \\ &+ \max_{1 \leq j \leq N} \hat{P}_n w(\cdot, \delta_{c_j}, c_j) + \max_{1 \leq j \leq N} P_\theta w(\cdot, \delta_{c_j}, c_j) \\ &\rightarrow 2 \max_{1 \leq j \leq N} P_\theta w(\cdot, \delta_{c_j}, c_j) \leq 2\epsilon, \mathbf{P}_\theta\text{-a.s.} \end{aligned}$$

□

Example The above theorem directly uses the definition of the M-estimator, and thus does not rely on having an explicit expression available. Here is an example where an explicit expression is indeed not possible. Consider the logistic location family, where the densities are

$$p_\theta(x) = \frac{e^{x-\theta}}{(1 + e^{x-\theta})^2}, \quad x \in \mathbb{R},$$

where $\theta \in \Theta \subset \mathbb{R}$ is the location parameter. Take

$$\rho_\theta(x) := -\log p_\theta(x) = \theta - x + 2 \log(1 + e^{x-\theta}).$$

So $\hat{\theta}_n$ is a solution of

$$\frac{2}{n} \sum_{i=1}^n \frac{e^{X_i - \hat{\theta}_n}}{1 + e^{X_i - \hat{\theta}_n}} = 1.$$

This expression cannot be made into an explicit expression. However, we do note the caveat that in order to be able to apply the above consistency theorem, we need to assume that Θ is bounded. This problem can be circumvented by using the result below for Z-estimators.

To prove consistency of a Z-estimator of a one-dimensional parameter is relatively easy.

Theorem 6.3.2 *Assume that $\Gamma \subset \mathbb{R}$, that $\psi_c(x)$ is continuous in c for all x , that*

$$P_\theta |\psi_c| < \infty, \quad \forall c,$$

and that $\exists \delta > 0$ such that

$$P_\theta \psi_c > 0, \quad \gamma < c < \gamma + \delta,$$

$$P_\theta \psi_c < 0, \quad \gamma - \delta < c < \gamma.$$

Then for n large enough, \mathbf{P}_θ -a.s., there is a solution $\hat{\gamma}_n$ of $\hat{P}_n \psi_{\hat{\gamma}_n} = 0$, and this solution $\hat{\gamma}_n$ is consistent.

Proof. Let $0 < \epsilon < \delta$ be arbitrary. By the law of large numbers, for n sufficiently large, \mathbf{P}_θ -a.s.,

$$\hat{P}_n \psi_{\gamma+\epsilon} > 0, \quad \hat{P}_n \psi_{\gamma-\epsilon} < 0.$$

The continuity of $c \mapsto \psi_c$ implies that then $\hat{P}_n \psi_{\hat{\gamma}_n} = 0$ for some $|\hat{\gamma}_n - \gamma| < \epsilon$. □

6.3.2 Asymptotic normality of M-estimators

Recall the CLT: for each $f : \mathcal{X} \rightarrow \mathbb{R}^r$ for which

$$\Sigma := P_\theta f f^T - (P_\theta f)(P_\theta f)^T$$

exists, we have

$$\sqrt{n}(\hat{P}_n - P_\theta)f \xrightarrow{D_\theta} \mathcal{N}(0, \Sigma).$$

Denote now

$$\nu_n(c) := \sqrt{n}(\hat{P}_n - P_\theta)\psi_c, \quad c \in \Gamma.$$

Definition *The stochastic process*

$$\{\nu_n(c) : c \in \Gamma\}$$

is called the empirical process indexed by c . The empirical process is called asymptotically continuous at γ if for all (possibly random) sequences $\{\gamma_n\}$ in Γ , with $\|\gamma_n - \gamma\| = o_{\mathbf{P}_\theta}(1)$, we have

$$|\nu_n(\gamma_n) - \nu_n(\gamma)| = o_{\mathbf{P}_\theta}(1).$$

For verifying asymptotic continuity, there are various tools, which involve complexity assumptions on the map $c \mapsto \psi_c$. This goes beyond the scope of these notes. Asymptotic linearity can also be established directly, under rather restrictive assumptions, see Theorem 6.3.4 below. But first, let us see what asymptotic continuity can bring us.

We assume that

$$M_\theta := \left. \frac{\partial}{\partial c^T} P_\theta \psi_c \right|_{c=\gamma}$$

exists. It is a $p \times p$ matrix. We require it to be of full rank, which amounts to assuming that γ , as a solution to $P_\theta \psi_\gamma = 0$, is well-identified.

Theorem 6.3.3 *Let $\hat{\gamma}_n$ be the Z-estimator of γ , and suppose that $\hat{\gamma}_n$ is a consistent estimator of γ , and that ν_n is asymptotically continuous at γ . Suppose moreover M_θ^{-1} exists, and also*

$$J_\theta := P_\theta \psi_\gamma \psi_\gamma^T.$$

Then $\hat{\gamma}_n$ is asymptotically linear, with influence function

$$l_\theta = -M_\theta^{-1} \psi_\gamma.$$

Hence

$$\sqrt{n}(\hat{\gamma}_n - \gamma) \xrightarrow{D_\theta} \mathcal{N}(0, V_\theta),$$

with

$$V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1}.$$

Proof. By definition,

$$\hat{P}_n \psi_{\hat{\gamma}_n} = 0, \quad P_\theta \psi_\gamma = 0.$$

So we have

$$\begin{aligned} 0 &= \hat{P}_n \psi_{\hat{\gamma}_n} = (\hat{P}_n - P_\theta) \psi_{\hat{\gamma}_n} + P_\theta \psi_{\hat{\gamma}_n} \\ &= (\hat{P}_n - P_\theta) \psi_{\hat{\gamma}_n} + P_\theta (\psi_{\hat{\gamma}_n} - \psi_\gamma) \\ &= (i) + (ii). \end{aligned}$$

For the first term, we use the asymptotic continuity of ν_n at γ :

$$\begin{aligned} (i) &= (\hat{P}_n - P_\theta) \psi_{\hat{\gamma}_n} = \nu_n(\hat{\gamma}_n) / \sqrt{n} = \nu_n(\gamma) / \sqrt{n} + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \\ &= \hat{P}_n \psi_\gamma + o_{\mathbf{P}_\theta}(1/n). \end{aligned}$$

For the second term, we use the differentiability of $P_\theta \psi_c$ at $c = \gamma$:

$$(ii) = P_\theta (\psi_{\hat{\gamma}_n} - \psi_\gamma) = M(\hat{\gamma}_n - \gamma) + o(\|\hat{\gamma}_n - \gamma\|).$$

So we arrive at

$$0 = \hat{P}_n \psi_\gamma + o_{\mathbf{P}_\theta}(1/n) + M(\hat{\gamma}_n - \gamma) + o(\|\hat{\gamma}_n - \gamma\|).$$

Because, by the CLT, $\hat{P}_n \psi_\gamma = O_{\mathbf{P}_\theta}(1/\sqrt{n})$, this implies $\|\hat{\gamma}_n - \gamma\| = O_{\mathbf{P}_\theta}(1/\sqrt{n})$. Hence

$$0 = \hat{P}_n \psi_\gamma + M(\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}),$$

or

$$M(\hat{\gamma}_n - \gamma) = -\hat{P}_n \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}),$$

or

$$(\hat{\gamma}_n - \gamma) = -\hat{P}_n M^{-1} \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}).$$

□

In the next theorem, we assume quite a lot of smoothness for the functions ψ_c (namely, derivatives that are Lipschitz), so that asymptotic linearity can be proved by straightforward arguments. We stress however that such smoothness assumptions are by no means necessary.

Theorem 6.3.4 *Let $\hat{\gamma}_n$ be the Z-estimator of γ , and suppose that $\hat{\gamma}_n$ is a consistent estimator of γ . Suppose that, for all c in a neighborhood $\{c \in \Gamma : \|c - \gamma\| < \epsilon\}$, the map $c \mapsto \psi_c(x)$ is differentiable for all x , with derivative*

$$\dot{\psi}_c(x) = \frac{\partial}{\partial c^T} \psi_c(x)$$

(a $p \times p$ matrix). Assume moreover that, for all c and \tilde{c} in a neighborhood of γ , and for all x , we have, in matrix-norm⁴,

$$\|\dot{\psi}_c(x) - \dot{\psi}_{\tilde{c}}(x)\| \leq H(x) \|c - \tilde{c}\|,$$

⁴For a matrix A , $\|A\| := \sup_{v \neq 0} \|Av\|/\|v\|$.

where $H : \mathcal{X} \rightarrow \mathbb{R}$ satisfies

$$P_\theta H < \infty.$$

Then

$$M_\theta = \frac{\partial}{\partial c^T} P_\theta \psi_c \Big|_{c=\gamma} = P_\theta \dot{\psi}_\gamma. \quad (6.4)$$

Assuming M^{-1} and $J := E_\theta \psi_\gamma \psi_\gamma^T$ exist, the influence function of $\hat{\gamma}_n$ is

$$l_\theta = -M_\theta^{-1} \psi_\gamma.$$

Proof. Result (6.4) follows from the dominated convergence theorem.

By the mean value theorem,

$$0 = \hat{P}_n \psi_{\hat{\gamma}_n} = \hat{P}_n \psi_\gamma + \hat{P}_n \dot{\psi}_{\tilde{\gamma}_n(\cdot)} (\hat{\gamma}_n - \gamma)$$

where for all x , $\|\tilde{\gamma}_n(x) - \gamma\| \leq \|\hat{\gamma}_n - \gamma\|$. Thus

$$0 = \hat{P}_n \psi_\gamma + \hat{P}_n \dot{\psi}_\gamma (\hat{\gamma}_n - \gamma) + \hat{P}_n (\dot{\psi}_{\tilde{\gamma}_n(\cdot)} - \dot{\psi}_\gamma) (\hat{\gamma}_n - \gamma),$$

so that

$$\left| \hat{P}_n \psi_\gamma + \hat{P}_n \dot{\psi}_\gamma (\hat{\gamma}_n - \gamma) \right| \leq \hat{P}_n H \|\hat{\gamma}_n - \gamma\|^2 = O_{\mathbf{P}_\theta}(1) \|\hat{\gamma}_n - \gamma\|^2,$$

where in the last inequality, we used $P_\theta H < \infty$. Now, by the law of large numbers,

$$\hat{P}_n \dot{\psi}_\gamma = P_\theta \dot{\psi}_\gamma + o_{\mathbf{P}_\theta}(1) = M_\theta + o_{\mathbf{P}_\theta}(1).$$

Thus

$$\left| \hat{P}_n \psi_\gamma + M_\theta (\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(\|\hat{\gamma}_n - \gamma\|) \right| = O_{\mathbf{P}_\theta}(\|\hat{\gamma}_n - \gamma\|^2).$$

Because $\hat{P}_n \psi_\gamma = O_{\mathbf{P}_\theta}(1/\sqrt{n})$, this ensures that $\|\hat{\gamma}_n - \gamma\| = O_{\mathbf{P}_\theta}(1/\sqrt{n})$. It follows that

$$\left| \hat{P}_n \psi_\gamma + M_\theta (\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \right| = O_{\mathbf{P}_\theta}(1/n).$$

Hence

$$M_\theta (\hat{\gamma}_n - \gamma) = -\hat{P}_n \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n})$$

and so

$$(\hat{\gamma}_n - \gamma) = -\hat{P}_n M_\theta^{-1} \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}).$$

□

Example 6.3.3 In this example, we show that, under regularity conditions, the MLE is asymptotically normal with asymptotic covariance matrix the inverse of the Fisher-information matrix $I(\theta)$. Let $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ be dominated by a σ -finite dominating measure ν , and write the densities as $p_\theta = dP_\theta/d\nu$.

Suppose that $\Theta \subset \mathbb{R}^p$. Assume condition I, i.e. that the support of p_θ does not depend on θ . As loss we take minus the log-likelihood:

$$\rho_\theta := -\log p_\theta.$$

We suppose that the score function

$$s_\theta = \frac{\partial}{\partial \theta} \log p_\theta = \frac{\dot{p}_\theta}{p_\theta}$$

exists, and that we may interchange differentiation and integration, so that the score has mean zero.

$$P_\theta s_\theta = \int \dot{p}_\theta d\nu = \frac{\partial}{\partial \theta} \int p_\theta d\nu = \frac{\partial}{\partial \theta} 1 = 0.$$

Recall that the Fisher-information matrix is

$$I(\theta) := P_\theta s_\theta s_\theta^T.$$

Now, it is clear that $\psi_\theta = -s_\theta$, and, assuming derivatives exist and that again we may change the order of differentiation and integration,

$$M_\theta = P_\theta \dot{\psi}_\theta = -P_\theta \dot{s}_\theta,$$

and

$$\begin{aligned} P_\theta \dot{s}_\theta &= P_\theta \left(\frac{\ddot{p}_\theta}{p_\theta} - s_\theta s_\theta^T \right) \\ &= \left(\frac{\partial^2}{\partial \theta \partial \theta^T} 1 \right) - P_\theta s_\theta s_\theta^T \\ &= 0 - I(\theta). \end{aligned}$$

Hence, in this case, $M_\theta = -I(\theta)$, and the influence function of the MLE

$$\hat{\theta}_n := \arg \max_{\hat{\theta} \in \Theta} \hat{P}_n \log p_{\hat{\theta}}$$

is

$$l_\theta = I(\theta)^{-1} s_\theta.$$

So the asymptotic covariance matrix of the MLE $\hat{\theta}_n$ is

$$I(\theta)^{-1} \left(P_\theta s_\theta s_\theta^T \right) I(\theta)^{-1} = I(\theta)^{-1}.$$

Example 6.3.4 In this example, the parameter of interest is the α -quantile. We will consider a loss function which does not satisfy regularity conditions, but nevertheless leads to an asymptotically linear estimator.



Let $\mathcal{X} := \mathbb{R}$. The distribution function of X is denoted by F . Let $0 < \alpha < 1$ be given. The α -quantile of F is $\gamma = F^{-1}(\alpha)$ (assumed to exist). We moreover

assume that F has density f with respect to Lebesgue measure, and that $f(x) > 0$ in a neighborhood of γ . As loss function we take

$$\rho_c(x) := \rho(x - c),$$

where

$$\rho(x) := (1 - \alpha)|x|1\{x < 0\} + \alpha|x|1\{x > 0\}.$$

We now first check that

$$\arg \min_c P_\theta \rho_c = F^{-1}(\alpha) := \gamma.$$

We have

$$\dot{\rho}(x) = \alpha 1\{x > 0\} - (1 - \alpha)1\{x < 0\}.$$

Note that $\dot{\rho}$ does not exist at $x = 0$. This is one of the irregularities in this example.

It follows that

$$\psi_c(x) = -\alpha 1\{x > c\} + (1 - \alpha)1\{x < c\}.$$

Hence

$$P_\theta \psi_c = -\alpha + F(c)$$

(the fact that ψ_c is not defined at $x = c$ can be shown not to be a problem, roughly because a single point has probability zero, as F is assumed to be continuous). So

$$P_\theta \psi_\gamma = 0, \text{ for } \gamma = F^{-1}(\alpha).$$

We now derive M_θ , which is a scalar in this case:

$$\begin{aligned} M_\theta &= \left. \frac{d}{dc} P_\theta \psi_c \right|_{c=\gamma} \\ &= \left. \frac{d}{dc} (-\alpha + F(c)) \right|_{c=\gamma} = f(\gamma) = f(F^{-1}(\alpha)). \end{aligned}$$

The influence function is thus ⁵

$$l_\theta(x) = -M_\theta^{-1} \psi_\gamma(x) = \frac{1}{f(\gamma)} \left\{ -1\{x < \gamma\} + \alpha \right\}.$$

We conclude that, for

$$\hat{\gamma}_n = \arg \min_c \hat{P}_n \rho_c,$$

which we write as the sample quantile $\hat{\gamma}_n = \hat{F}_n^{-1}(\alpha)$ (or an approximation thereof up to order $o_{\mathbf{P}_\theta}(1/\sqrt{n})$), one has

$$\sqrt{n}(\hat{F}_n^{-1}(\alpha) - F^{-1}(\alpha)) \xrightarrow{D_\theta} \mathcal{N}\left(0, \frac{\alpha(1 - \alpha)}{f^2(F^{-1}(\alpha))}\right).$$

⁵Note that in the special case $\alpha = 1/2$ (where γ is the median), this becomes

$$l_\theta(x) = \begin{cases} -\frac{1}{2f(\gamma)} & x < \gamma \\ +\frac{1}{2f(\gamma)} & x > \gamma \end{cases}.$$

Example 6.3.5 In this example, we illustrate that the Huber-estimator is asymptotically linear. Let again $\mathcal{X} = \mathbb{R}$ and F be the distribution function of X . We let the parameter of interest be the a location parameter. The Huber loss function is

$$\rho_c(x) = \rho(x - c),$$

with

$$\rho(x) = \begin{cases} x^2 & |x| \leq k \\ k(2|x| - k) & |x| > k \end{cases}.$$

We define γ as

$$\gamma := \arg \min_c P_\theta \rho_c.$$

It holds that

$$\dot{\rho}(x) = \begin{cases} 2x & |x| \leq k \\ +2k & x > k \\ -2k & x < -k \end{cases}.$$

Therefore,

$$\psi_c(x) = \begin{cases} -2(x - c) & |x - c| \leq k \\ -2k & x - c > k \\ +2k & x - c < -k \end{cases}.$$

One easily derives that

$$\begin{aligned} P_\theta \psi_c &= -2 \int_{-k+c}^{k+c} x dF(x) + 2c[F(k+c) - F(-k+c)] \\ &\quad - 2k[1 - F(k+c)] + 2kF(-k+c). \end{aligned}$$

So

$$M_\theta = \frac{d}{dc} P_\theta \psi_c \Big|_{c=\gamma} = 2[F(k+\gamma) - F(-k+\gamma)].$$

The influence function of the Huber estimator is

$$l_\theta(x) = \frac{1}{[F(k+\gamma) - F(-k+\gamma)]} \begin{cases} x - \gamma & |x - \gamma| \leq k \\ +k & x - \gamma > k \\ -k & x - \gamma < -k \end{cases}.$$

For $k \rightarrow 0$, this corresponds to the influence function of the median.

6.4 Plug-in estimators

When \mathcal{X} is Euclidean space, one can define the distribution function $F(x) := P_\theta(X \leq x)$ and the empirical distribution function

$$\hat{F}_n(x) = \frac{1}{n} \#\{X_i \leq x, 1 \leq i \leq n\}.$$

This is the distribution function of a probability measure that puts mass $1/n$ at each observation. For general \mathcal{X} , we define likewise the empirical distribution \hat{P}_n as the distribution that puts mass $1/n$ at each observation, i.e., more formally

$$\hat{P}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i},$$

where δ_x is a point mass at x . Thus, for (measurable) sets $A \subset \mathcal{X}$,

$$\hat{P}_n(A) = \frac{1}{n} \#\{X_i \in A, 1 \leq i \leq n\}.$$

For (measurable) functions $f : \mathcal{X} \rightarrow \mathbb{R}^r$, we write, as in the previous section,

$$\hat{P}_n f := \frac{1}{n} \sum_{i=1}^n f(X_i) = \int f d\hat{P}_n.$$

Thus, for sets,

$$\hat{P}_n(A) = \hat{P}_n 1_A.$$

Again, as in the previous section, we use the same notations for expectations under P_θ :

$$P_\theta f := E_\theta f(X) = \int f dP_\theta,$$

so that

$$P_\theta(A) = P_\theta 1_A.$$

The parameter of interest is denoted as

$$\gamma = g(\theta) \in \mathbb{R}^p.$$

It can often be written in the form

$$\gamma = Q(P_\theta),$$

where Q is some functional on (a subset of) the model class \mathcal{P} . Assuming Q is also defined at the empirical measure \hat{P}_n , the plug-in estimator of γ is now

$$T_n := Q(\hat{P}_n).$$

Conversely,

Definition *If a statistic T_n can be written as $T_n = Q(\hat{P}_n)$, then it is called a Fisher-consistent estimator of $\gamma = g(\theta)$, if $Q(P_\theta) = g(\theta)$ for all $\theta \in \Theta$.*

We will also encounter modifications, where

$$T_n = Q_n(\hat{P}_n),$$

and for n large,

$$Q_n(P_\theta) \approx Q(P_\theta) = g(\theta).$$

Example Let $\gamma := h(P_\theta f)$. The plug-in estimator is then $T_n = h(\hat{P}_n f)$.

Example The M-estimator $\hat{\gamma}_n = \arg \min_{c \in \Gamma} \hat{P}_n \rho_c$ is a plug-in estimator of $\gamma = \arg \min_{c \in \Gamma} P_\theta \rho_c$ (and similarly for the Z-estimator).

Example Let $\mathcal{X} = \mathbb{R}$ and consider the α -trimmed mean

$$T_n := \frac{1}{n - 2[n\alpha]} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} X_{(i)}.$$

What is its theoretical counterpart? Because the i -th order statistic $X_{(i)}$ can be written as

$$X_{(i)} = \hat{F}_n^{-1}(i/n),$$

and in fact

$$X_{(i)} = \hat{F}_n^{-1}(u), \quad i/n \leq u < (i+1)/n,$$

we may write, for $\alpha_n := [n\alpha]/n$,

$$\begin{aligned} T_n &= \frac{n}{n - 2[n\alpha]} \frac{1}{n} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} \hat{F}_n^{-1}(i/n) \\ &= \frac{1}{1 - 2\alpha_n} \int_{\alpha_n+1/n}^{1-\alpha_n} \hat{F}_n^{-1}(u) du := Q_n(\hat{P}_n). \end{aligned}$$

Replacing \hat{F}_n by F gives,

$$\begin{aligned} Q_n(F) &= \frac{1}{1 - 2\alpha_n} \int_{\alpha_n+1/n}^{1-\alpha_n} F^{-1}(u) du \\ &\approx \frac{1}{1 - 2\alpha} \int_{\alpha}^{1-\alpha} F^{-1}(u) du = \frac{1}{1 - 2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x) := Q(P_\theta). \end{aligned}$$

Example Let $\mathcal{X} = \mathbb{R}$, and suppose X has density f w.r.t., Lebesgue measure. Suppose f is the parameter of interest. We may write

$$f(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h}.$$

Replacing F by \hat{F}_n here does not make sense. Thus, this is an example where $Q(P) = f$ is only well defined for distributions P that have a density f . We may however slightly extend the plug-in idea, by using the estimator

$$\hat{f}_n(x) = \frac{\hat{F}_n(x+h_n) - \hat{F}_n(x-h_n)}{2h_n} := Q_n(\hat{P}_n),$$

with h_n “small” ($h_n \rightarrow 0$ as $n \rightarrow \infty$).

6.4.1 Consistency of plug-in estimators

We first present the uniform convergence of the empirical distribution function to the theoretical one.

Such uniform convergence results hold also in much more general settings (see also (6.3) in the proof of consistency for M-estimators).

Theorem 6.4.1 (*Glivenko-Cantelli*) *Let $\mathcal{X} = \mathbb{R}$. We have*

$$\sup_x |\hat{F}_n(x) - F(x)| \rightarrow 0, \mathbb{P}_\theta\text{-a.s..}$$

Proof. We know that by the law of large numbers, for all x

$$|\hat{F}_n(x) - F(x)| \rightarrow 0, \mathbb{P}_\theta\text{-a.s.,}$$

so also for all finite collection a_1, \dots, a_N ,

$$\max_{1 \leq j \leq N} |\hat{F}_n(a_j) - F(a_j)| \rightarrow 0, \mathbb{P}_\theta\text{-a.s..}$$

Let $\epsilon > 0$ be arbitrary, and take $a_0 < a_1 < \dots < a_{N-1} < a_N$ in such a way that

$$F(a_j) - F(a_{j-1}) \leq \epsilon, j = 1, \dots, N$$

where $F(a_0) := 0$ and $F(a_N) := 1$. Then, when $x \in (a_{j-1}, a_j]$,

$$\hat{F}_n(x) - F(x) \leq \hat{F}_n(a_j) - F(a_{j-1}) \leq F_n(a_j) - F(a_j) + \epsilon,$$

and

$$\hat{F}_n(x) - F(x) \geq \hat{F}_n(a_{j-1}) - F(a_j) \geq \hat{F}_n(a_{j-1}) - F(a_{j-1}) - \epsilon,$$

so

$$\sup_x |\hat{F}_n(x) - F(x)| \leq \max_{1 \leq j \leq N} |\hat{F}_n(a_j) - F(a_j)| + \epsilon \rightarrow \epsilon, \mathbb{P}_\theta\text{-a.s..}$$

□

Example Let $\mathcal{X} = \mathbb{R}$ and let F be the distribution function of X . We consider estimating the median $\gamma := F^{-1}(1/2)$. We assume F to continuous and strictly increasing. The sample median is

$$T_n := \hat{F}_n^{-1}(1/2) := \begin{cases} X_{((n+1)/2)} & n \text{ odd} \\ [X_{(n/2)} + X_{(n/2+1)}]/2 & n \text{ even} \end{cases}.$$

So

$$\hat{F}_n(T_n) = \frac{1}{2} + \begin{cases} 1/(2n) & n \text{ odd} \\ 0 & n \text{ even} \end{cases}.$$

It follows that

$$\begin{aligned} |F(T_n) - F(\gamma)| &\leq |\hat{F}_n(T_n) - F(T_n)| + |\hat{F}_n(T_n) - F(\gamma)| \\ &= |\hat{F}_n(T_n) - F(T_n)| + |\hat{F}_n(T_n) - \frac{1}{2}| \\ &\leq |\hat{F}_n(T_n) - F(T_n)| + \frac{1}{2n} \rightarrow 0, \mathbb{P}_\theta\text{-a.s..} \end{aligned}$$

So $\hat{F}_n^{-1}(1/2) = T_n \rightarrow \gamma = F^{-1}(1/2)$, $\mathbb{P}_\theta\text{-a.s.}$, i.e., the sample median is a consistent estimator of the population median.

6.4.2 Asymptotic normality of plug-in estimators

Let $\gamma := Q(P) \in \mathbb{R}^p$ be the parameter of interest. The idea in this subsection is to apply a δ -method, but now in a nonparametric framework. The parametric δ -method says that if $\hat{\theta}_n$ is an asymptotically linear estimator of $\theta \in \mathbb{R}^p$, and if $\gamma = g(\theta)$ is some function of the parameter θ , with g being differentiable at θ , then $\hat{\gamma}$ is an asymptotically linear estimator of γ . Now, we write $\gamma = Q(P)$ as a function of the probability measure P (with $P = P_\theta$, so that $g(\theta) = Q(P_\theta)$). We let P play the role of θ , i.e., we use the probability measures themselves as parameterization of \mathcal{P} . We then have to redefine differentiability in an abstract setting, namely we differentiate w.r.t. P .

Definition

◦ The influence function of Q at P is

$$l_P(x) := \lim_{\epsilon \downarrow 0} \frac{Q((1-\epsilon)P + \epsilon\delta_x) - Q(P)}{\epsilon}, \quad x \in \mathcal{X},$$

whenever the limit exists.

◦ The map Q is called Gâteaux differentiable at P if for all probability measures \tilde{P} , we have

$$\lim_{\epsilon \downarrow 0} \frac{Q((1-\epsilon)P + \epsilon\tilde{P}) - Q(P)}{\epsilon} = E_{\tilde{P}}l_P(X).$$

◦ Let d be some (pseudo-)metric on the space of probability measures. The map Q is called Fréchet differentiable at P , with respect to the metric d , if

$$Q(\tilde{P}) - Q(P) = E_{\tilde{P}}l_P(X) + o(d(\tilde{P}, P)).$$

Remark 1 In line with the notation introduced previously, we write for a function $f : \mathcal{X} \rightarrow \mathbb{R}^r$ and a probability measure \tilde{P} on \mathcal{X}

$$\tilde{P}f := E_{\tilde{P}}f(X).$$

Remark 2 If Q is Fréchet or Gâteaux differentiable at P , then

$$Pl_P(:= E_Pl_P(X)) = 0.$$

Remark 3 If Q is Fréchet differentiable at P , and if moreover

$$d((1-\epsilon)P + \epsilon\tilde{P}, P) = o(\epsilon), \quad \epsilon \downarrow 0,$$

then Q is Gâteaux differentiable at P :

$$\begin{aligned} Q((1-\epsilon)P + \epsilon\tilde{P}) - Q(P) &= ((1-\epsilon)P + \epsilon\tilde{P})l_P + o(\epsilon) \\ &= \epsilon\tilde{P}l_P + o(\epsilon). \end{aligned}$$

We now show that Fréchet differentiable functionals are generally asymptotically linear.

Lemma 6.4.1 *Suppose that Q is Fréchet differentiable at P with influence function l_P , and that*

$$d(\hat{P}_n, P) = O_{\mathbf{P}}(n^{-1/2}). \quad (6.5)$$

Then

$$Q(\hat{P}_n) - Q(P) = \hat{P}_n l_P + o_{\mathbf{P}}(n^{-1/2}).$$

Proof. This follows immediately from the definition of Fréchet differentiability. \square

Corollary 6.4.1 *Assume the conditions of Lemma 6.4.1, with influence function l_P satisfying $V_P := Pl_P l_P^T < \infty$. Then*

$$\sqrt{n}(Q(\hat{P}_n) - Q(P)) \xrightarrow{D_P} \mathcal{N}(0, V_P).$$

An example where (6.5) holds

Suppose $\mathcal{X} = \mathbb{R}$ and that we take

$$d(\tilde{P}, P) := \sup_x |\tilde{F}(x) - F(x)|.$$

Then indeed $d(\hat{P}_n, P) = O_{\mathbf{P}}(n^{-1/2})$. This follows from Donsker's theorem, which we state here without proof:

Donsker's theorem *Suppose F is continuous. Then*

$$\sup_x \sqrt{n}|\hat{F}_n(x) - F(x)| \xrightarrow{D} Z,$$

where the random variable Z has distribution function

$$G(z) = 1 - 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp[-2j^2 z^2], \quad z \geq 0.$$

Fréchet differentiability is generally quite hard to prove, and often not even true. We will only illustrate Gâteaux differentiability in some examples.

Example 6.4.1 We consider the Z-estimator. Throughout in this example, we assume enough regularity.



Let γ be defined by the equation

$$P\psi_{\gamma} = 0.$$

Let $P_{\epsilon} := (1 - \epsilon)P + \epsilon\tilde{P}$, and let γ_{ϵ} be a solution of the equation

$$P_{\epsilon}\psi_{\gamma_{\epsilon}} = 0.$$

We assume that as $\epsilon \downarrow 0$, also $\gamma_\epsilon \rightarrow \gamma$. It holds that

$$(1 - \epsilon)P\psi_{\gamma_\epsilon} + \epsilon\tilde{P}\psi_{\gamma_\epsilon} = 0,$$

so

$$P\psi_{\gamma_\epsilon} + \epsilon(\tilde{P} - P)\psi_{\gamma_\epsilon} = 0,$$

and hence

$$P(\psi_{\gamma_\epsilon} - \psi_\gamma) + \epsilon(\tilde{P} - P)\psi_{\gamma_\epsilon} = 0.$$

Assuming differentiability of $c \mapsto P\psi_c$, we obtain

$$\begin{aligned} P(\psi_{\gamma_\epsilon} - \psi_\gamma) &= \left(\frac{\partial}{\partial c^T} P\psi_c \Big|_{c=\gamma} \right) (\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|) \\ &:= M_P(\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|). \end{aligned}$$

Moreover, again under regularity

$$\begin{aligned} (\tilde{P} - P)\psi_{\gamma_\epsilon} &= (\tilde{P} - P)\psi_\gamma + (\tilde{P} - P)(\psi_{\gamma_\epsilon} - \psi_\gamma) \\ &= (\tilde{P} - P)\psi_\gamma + o(1) = \tilde{P}\psi_\gamma + o(1). \end{aligned}$$

It follows that

$$M_P(\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|) + \epsilon(\tilde{P} - P)\psi_\gamma + o(\epsilon) = 0,$$

or, assuming M_P to be invertible,

$$(\gamma_\epsilon - \gamma)(1 + o(1)) = -\epsilon M_P^{-1} \tilde{P}\psi_\gamma + o(\epsilon),$$

which gives

$$\frac{\gamma_\epsilon - \gamma}{\epsilon} \rightarrow M_P^{-1} \tilde{P}\psi_\gamma.$$

The influence function is thus (as already seen in Subsection 6.3.2)

$$l_P = M_P^{-1} \tilde{P}\psi_\gamma.$$

Example 6.4.2 The α -trimmed mean is a plug-in estimator of

$$\gamma := Q(P) = \frac{1}{1 - 2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x).$$

Using partial integration, may write this as

$$(1 - 2\alpha)\gamma = (1 - \alpha)F^{-1}(1 - \alpha) - \alpha F^{-1}(\alpha) - \int_\alpha^{1-\alpha} v dF^{-1}(v).$$

The influence function of the quantile $F^{-1}(v)$ is

$$q_v(x) = -\frac{1}{f(F^{-1}(v))} \left(1\{x \leq F^{-1}(v)\} - v \right)$$

(see Example 6.3.4), i.e., for the distribution $P_\epsilon = (1 - \epsilon)P + \epsilon\tilde{P}$, with distribution function $F_\epsilon = (1 - \epsilon)F + \epsilon\tilde{F}$, we have

$$\lim_{\epsilon \downarrow 0} \frac{F_\epsilon^{-1}(v) - F^{-1}(v)}{\epsilon} = \tilde{P}q_v = -\frac{1}{f(F^{-1}(v))} \left(\tilde{F}(F^{-1}(v)) - v \right).$$

Hence, for $P_\epsilon = (1 - \epsilon)P + \epsilon\tilde{P}$,

$$\begin{aligned} (1 - 2\alpha) \lim_{\epsilon \downarrow 0} \frac{Q((1 - \epsilon)P + \epsilon\tilde{P}) - Q(P)}{\epsilon} &= (1 - \alpha)\tilde{P}q_{1-\alpha} - \alpha\tilde{P}q_\alpha - \int_\alpha^{1-\alpha} v d\tilde{P}q_v \\ &= \int_\alpha^{1-\alpha} \frac{1}{f(F^{-1}(v))} \left(\tilde{F}(F^{-1}(v)) - v \right) dv \\ &= \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \frac{1}{f(u)} \left(\tilde{F}(u) - F(u) \right) dF(u) = \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \left(\tilde{F}(u) - F(u) \right) du \\ &= (1 - 2\alpha)\tilde{P}l_P, \end{aligned}$$

where

$$l_P(x) = -\frac{1}{1 - 2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \left(1\{x \leq u\} - F(u) \right) du.$$

We conclude that, under regularity conditions, the α -trimmed mean is asymptotically linear with the above influence function l_P , and hence asymptotically normal with asymptotic variance Pl_P^2 .

6.5 Asymptotic relative efficiency

In this section, we assume that the parameter of interest is real-valued:

$$\gamma \in \Gamma \subset \mathbb{R}.$$

Definition Let $T_{n,1}$ and $T_{n,2}$ be two estimators of γ , that satisfy

$$\sqrt{n}(T_{n,j} - \gamma) \xrightarrow{D_\theta} \mathcal{N}(0, V_{\theta,j}), \quad j = 1, 2.$$

Then

$$e_{2:1} := \frac{V_{\theta,1}}{V_{\theta,2}}$$

is called the asymptotic relative efficiency of $T_{n,2}$ with respect to $T_{n,1}$.

If $e_{2:1} > 1$, the estimator $T_{n,2}$ is asymptotically more efficient than $T_{n,1}$. An asymptotic $(1 - \alpha)$ -confidence interval for γ based on $T_{n,2}$ is then narrower than the one based on $T_{n,1}$.

Example 6.5.1 Let $\mathcal{X} = \mathbb{R}$, and F be the distribution function of X . Suppose that F is symmetric around the parameter of interest μ . In other words,

$$F(\cdot) = F_0(\cdot - \mu),$$

where F_0 is symmetric around zero. We assume that F_0 has finite variance σ^2 , and that it has density f_0 w.r.t. Lebesgue measure, with $f_0(0) > 0$. Take $T_{n,1} := \bar{X}_n$, the sample mean, and $T_{n,2} := \hat{F}_n^{-1}(1/2)$, the sample median. Then $V_{\theta,1} = \sigma^2$ and $V_{\theta,2} = 1/(4f_0^2(0))$ (the latter being derived in Example 6.3.4). So

$$e_{2:1} = 4\sigma^2 f_0^2(0).$$

Whether the sample mean is the winner, or rather the sample median, depends thus on the distribution F_0 . Let us consider three cases.

Case i Let F_0 be the standard normal distribution, i.e., $F_0 = \Phi$. Then $\sigma^2 = 1$ and $f_0(0) = 1/\sqrt{2\pi}$. Hence

$$e_{2:1} = \frac{2}{\pi} \approx 0.64.$$

So \bar{X}_n is the winner. Note that \bar{X}_n is the MLE in this case.

Case ii Let F_0 be the Laplace distribution, with variance σ^2 equal to one. This distribution has density

$$f_0(x) = \frac{1}{\sqrt{2}} \exp[-\sqrt{2}|x|], \quad x \in \mathbb{R}.$$

So we have $f_0(0) = 1/\sqrt{2}$, and hence

$$e_{2:1} = 2.$$

Thus, the sample median, which is the MLE for this case, is the winner.

Case iii Suppose

$$F_0 = (1 - \eta)\Phi + \eta\Phi(\cdot/3).$$

This means that the distribution of X is a mixture, with proportions $1 - \eta$ and η , of two normal distributions, one with unit variance, and one with variance 3^2 . Otherwise put, associated with X is an unobservable label $Y \in \{0, 1\}$. If $Y = 1$, the random variable X is $\mathcal{N}(\mu, 1)$ -distributed. If $Y = 0$, the random variable X has a $\mathcal{N}(\mu, 3^2)$ distribution. Moreover, $P(Y = 1) = 1 - P(Y = 0) = 1 - \eta$. Hence

$$\sigma^2 := \text{var}(X) = (1 - \eta)\text{var}(X|Y = 1) + \eta\text{var}(X|Y = 0) = (1 - \eta) + 9\eta = 1 + 8\eta.$$

It furthermore holds that

$$f_0(0) = (1 - \eta)\phi(0) + \frac{\eta}{3}\phi(0) = \frac{1}{\sqrt{2\pi}} \left(1 - \frac{2\eta}{3}\right).$$

It follows that

$$e_{2:1} = \frac{2}{\pi} \left(1 - \frac{2\eta}{3}\right)^2 (1 + 8\eta).$$

Let us now further compare the results with the α -trimmed mean. Because F is symmetric, the α -trimmed mean has the same influence function as the Huber-estimator with $k = F^{-1}(1 - \alpha)$:

$$l_{\theta}(x) = \frac{1}{F_0(k) - F_0(-k)} \begin{cases} x - \mu, & |x - \mu| \leq k \\ +k, & x - \mu > k \\ -k, & x - \mu < -k \end{cases} .$$

This can be seen from Example 6.4.2. The influence function is used to compute the asymptotic variance $V_{\theta,\alpha}$ of the α -trimmed mean:

$$V_{\theta,\alpha} = \frac{\int_{F_0^{-1}(\alpha)}^{F_0^{-1}(1-\alpha)} x^2 dF_0(x) + 2\alpha(F_0^{-1}(1 - \alpha))^2}{(1 - 2\alpha)^2} .$$

From this, we then calculate the asymptotic relative efficiency of the α -trimmed mean w.r.t. the mean. Note that the median is the limiting case with $\alpha \rightarrow 1/2$.

Table: Asymptotic relative efficiency of α -trimmed mean over mean

	$\alpha = 0.05$	0.125	0.5
$\eta = 0.00$	0.99	0.94	0.64
0.05	1.20	1.19	0.83
0.25	1.40	1.66	1.33

6.6 Asymptotic Cramer Rao lower bound

Let X have distribution $P \in \{P_{\theta} : \theta \in \Theta\}$. We assume for simplicity that $\Theta \subset \mathbb{R}$ and that θ is the parameter of interest. Let T_n be an estimator of θ .

Throughout this section, we take certain, sometimes unspecified, regularity conditions for granted.



In particular, we assume that \mathcal{P} is dominated by some σ -finite measure ν , and that the Fisher-information

$$I(\theta) := E_{\theta} s_{\theta}^2(X)$$

exists for all θ . Here, s_{θ} is the score function

$$s_{\theta} := \frac{d}{d\theta} \log p_{\theta} = \dot{p}_{\theta}/p_{\theta},$$

with $p_{\theta} := dP_{\theta}/d\nu$.

Recall now that if T_n is an unbiased estimator of θ , then by the Cramer Rao lower bound, $1/I(\theta)$ is a lower bound for its variance (under regularity conditions I and II, see Section 3.3.1).

Definition Suppose that

$$\sqrt{n}(T_n - \theta) \xrightarrow{D_\theta} \mathcal{N}(b_\theta, V_\theta), \quad \forall \theta.$$

Then b_θ is called the asymptotic bias, and V_θ the asymptotic variance. The estimator T_n is called asymptotically unbiased if $b_\theta = 0$ for all θ . If T_n is asymptotically unbiased and moreover $V_\theta = 1/I(\theta)$ for all θ , and some regularity conditions holds, then T_n is called asymptotically efficient.

Remark 1 The assumptions in the above definition, are **for all** θ . Clearly, if one only looks at one fixed given θ_0 , it is easy to construct a super-efficient estimator, namely $T_n = \theta_0$. More generally, to avoid this kind of super-efficiency, one does not only require conditions to hold **for all** θ , but in fact **uniformly in** θ , or for all **sequences** $\{\theta_n\}$. The regularity one needs here involves the idea that one actually needs to allow for sequences θ_n the form $\theta_n = \theta + h/\sqrt{n}$. In fact, the regularity requirement is that also, for all h ,

$$\sqrt{n}(T_n - \theta_n) \xrightarrow{D_{\theta_n}} \mathcal{N}(0, V_\theta).$$

To make all this mathematically precise is quite involved. We refer to van der Vaart (1998). A glimpse is given in Le Cam's 3rd Lemma, see the next subsection.

Remark 2 Note that when $\theta = \theta_n$ is allowed to change with n , this means that distribution of X_i can change with n , and hence X_i can change with n . Instead of regarding the sample X_1, \dots, X_n are the first n of an infinite sequence, we now consider for each n a new sample, say $X_{1,1}, \dots, X_{n,n}$.

Remark 3 We have seen that the MLE $\hat{\theta}_n$ generally is indeed asymptotically unbiased with asymptotic variance V_θ equal to $1/I(\theta)$, i.e., under regularity assumptions, the MLE is asymptotically efficient.

For asymptotically linear estimators, with influence function l_θ , one has asymptotic variance $V_\theta = E_\theta l_\theta^2(X)$. The next lemma indicates that generally $1/I(\theta)$ is indeed a lower bound for the asymptotic variance.

Lemma 6.6.1 Suppose that

$$(T_n - \theta) = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbf{P}_\theta}(n^{-1/2}),$$

where $E_\theta l_\theta(X) = 0$, $E_\theta l_\theta^2(X) := V_\theta < \infty$. Assume moreover that

$$E_\theta l_\theta(X) s_\theta(X) = 1. \tag{6.6}$$

Then

$$V_\theta \geq \frac{1}{I(\theta)}.$$

Proof. This follows from the Cauchy-Schwarz inequality:

$$1 = |\text{cov}_\theta(l_\theta(X), s_\theta(X))|$$

$$\leq \text{var}_\theta(l_\theta(X))\text{var}_\theta(s_\theta(X)) = V_\theta I(\theta).$$

□

It may look like a coincidence when in a special case, equality (6.6) indeed holds. But actually, it is true in quite a few cases. This may at first seem like magic.

We consider two examples. To simplify the expressions, we again write shorthand

$$P_\theta f := E_\theta f(X).$$

Example 6.6.1 This example examines the Z-estimator of θ . Then we have, for $P = P_\theta$,

$$P\psi_\theta = 0.$$

The influence function is

$$l_\theta = -\psi_\theta/M_\theta,$$

where

$$M_\theta := \frac{d}{d\theta} P\psi_\theta.$$

Under regularity, we have

$$M_\theta = P\dot{\psi}_\theta = \int \dot{\psi}_\theta p_\theta d\nu, \quad \dot{\psi}_\theta = \frac{d}{d\theta} \psi_\theta.$$

We may also write

$$M_\theta = - \int \psi_\theta \dot{p}_\theta d\nu, \quad \dot{p}_\theta = \frac{d}{d\theta} p_\theta.$$

This follows from the chain rule

$$\frac{d}{d\theta} \psi_\theta p_\theta = \dot{\psi}_\theta p_\theta + \psi_\theta \dot{p}_\theta,$$

and (under regularity)

$$\int \frac{d}{d\theta} \psi_\theta p_\theta d\nu = \frac{d}{d\theta} \int \psi_\theta p_\theta d\nu = \frac{d}{d\theta} P\psi_\theta = \frac{d}{d\theta} 0 = 0.$$

Thus

$$Pl_\theta s_\theta = -M_\theta^{-1} P\psi_\theta s_\theta = -M_\theta^{-1} \int \psi_\theta \dot{p}_\theta d\nu = 1,$$

that is, (6.6) holds.

Example 6.6.2 We consider now the plug-in estimator $Q(\hat{P}_n)$. Suppose that Q is Fisher consistent (i.e., $Q(P_\theta) = \theta$ for all θ). Assume moreover that Q is Fréchet differentiable with respect to the metric d , at all P_θ , and that

$$d(P_{\tilde{\theta}}, P_\theta) = O(|\tilde{\theta} - \theta|).$$

Then, by the definition of Fréchet differentiability

$$h = Q(P_{\theta+h}) - Q(P_\theta) = P_{\theta+h} l_\theta + o(|h|) = (P_{\theta+h} - P_\theta) l_\theta + o(|h|),$$

or, as $h \rightarrow 0$,

$$\begin{aligned} 1 &= \frac{(P_{\theta+h} - P_{\theta})l_{\theta}}{h} + o(1) = \frac{\int l_{\theta}(p_{\theta+h} - p_{\theta})d\nu}{h} + o(1) \\ &\rightarrow \int l_{\theta}\dot{p}_{\theta}d\nu = P_{\theta}(l_{\theta}s_{\theta}). \end{aligned}$$

So (6.6) holds.

6.6.1 Le Cam's 3rd Lemma

The following example serves as a motivation to consider sequences θ_n depending on n . It shows that pointwise asymptotics can be very misleading.

Example 6.6.3 (*Hodges-Lehmann example of super-efficiency*) Let X_1, \dots, X_n be i.i.d. copies of X , where $X = \theta + \epsilon$, and ϵ is $\mathcal{N}(0, 1)$ -distributed. Consider the estimator

$$T_n := \begin{cases} \bar{X}_n, & \text{if } |\bar{X}_n| > n^{-1/4} \\ \bar{X}_n/2, & \text{if } |\bar{X}_n| \leq n^{-1/4} \end{cases}.$$

Then

$$\sqrt{n}(T_n - \theta) \xrightarrow{D_{\theta}} \begin{cases} \mathcal{N}(0, 1), & \theta \neq 0 \\ \mathcal{N}(0, \frac{1}{4}), & \theta = 0 \end{cases}.$$

So the pointwise asymptotics show that T_n can be more efficient than the sample average \bar{X}_n . But what happens if we consider sequences θ_n ? For example, let $\theta_n = h/\sqrt{n}$. Then, under \mathbb{P}_{θ_n} , $\bar{X}_n = \bar{\epsilon}_n + h/(\sqrt{n}) = O_{\mathbb{P}_{\theta_n}}(n^{-1/2})$. Hence, $\mathbb{P}_{\theta_n}(|\bar{X}_n| > n^{-1/4}) \rightarrow 0$, so that $\mathbb{P}_{\theta_n}(T_n = \bar{X}_n) \rightarrow 0$. Thus,

$$\begin{aligned} \sqrt{n}(T_n - \theta_n) &= \sqrt{n}(T_n - \theta_n)\mathbb{1}\{T_n = \bar{X}_n\} + \sqrt{n}(T_n - \theta_n)\mathbb{1}\{T_n = \bar{X}_n/2\} \\ &\xrightarrow{D_{\theta_n}} \mathcal{N}\left(-\frac{h}{2}, \frac{1}{4}\right). \end{aligned}$$

The asymptotic mean square error $\text{AMSE}_{\theta}(T_n)$ is defined as the asymptotic variance + asymptotic squared bias:

$$\text{AMSE}_{\theta_n}(T_n) = \frac{1 + h^2}{4}.$$

The $\text{AMSE}_{\theta}(\bar{X}_n)$ of \bar{X}_n is its normalized non-asymptotic mean square error, which is

$$\text{AMSE}_{\theta_n}(\bar{X}_n) = \text{MSE}_{\theta_n}(\bar{X}_n) = 1.$$

So when h is large enough, the asymptotic mean square error of T_n is larger than that of \bar{X}_n .

Le Cam's 3rd lemma shows that asymptotic linearity for all θ implies asymptotic normality, now also for sequences $\theta_n = \theta + h/\sqrt{n}$. The asymptotic variance for such sequences θ_n does not change. Moreover, if (6.6) holds for all θ , the estimator is also asymptotically unbiased under \mathbb{P}_{θ_n} .

Lemma 6.6.2 (*Le Cam's 3rd Lemma*) Suppose that for all θ ,

$$T_n - \theta = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbf{P}_\theta}(n^{-1/2}),$$

where $P_\theta l_\theta = 0$, and $V_\theta := P_\theta l_\theta^2 < \infty$. Then, under regularity conditions,

$$\sqrt{n}(T_n - \theta_n) \xrightarrow{D_{\theta_n}} \mathcal{N}\left(\{P_\theta(l_\theta s_\theta) - 1\}h, V_\theta\right).$$

We will present a sketch of the proof of this lemma. For this purpose, we need the following auxiliary lemma.

Lemma 6.6.3 (*Auxiliary lemma*) Let $Z \in \mathbb{R}^2$ be $\mathcal{N}(\mu, \Sigma)$ -distributed, where

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{1,2} \\ \sigma_{1,2} & \sigma_2^2 \end{pmatrix}.$$

Suppose that

$$\mu_2 = -\sigma_2^2/2.$$

Let $Y \in \mathbb{R}^2$ be $\mathcal{N}(\mu + a, \Sigma)$ -distributed, with

$$a = \begin{pmatrix} \sigma_{1,2} \\ \sigma_2^2 \end{pmatrix}.$$

Let ϕ_Z be the density of Z and ϕ_Y be the density of Y . Then we have the following equality for all $z = (z_1, z_2) \in \mathbb{R}^2$:

$$\phi_Z(z)e^{z_2^2} = \phi_Y(z).$$

Proof. The density of Z is

$$\phi_Z(z) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left[-\frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu)\right].$$

Now, one easily sees that

$$\Sigma^{-1}a = \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

So

$$\begin{aligned} \frac{1}{2}(z - \mu)^T \Sigma^{-1}(z - \mu) &= \frac{1}{2}(z - \mu - a)^T \Sigma^{-1}(z - \mu - a) \\ &\quad + a^T \Sigma^{-1}(z - \mu) - \frac{1}{2}a^T \Sigma^{-1}a \end{aligned}$$

and

$$\begin{aligned} a^T \Sigma^{-1}(z - \mu) - \frac{1}{2}a^T \Sigma^{-1}a &= \begin{pmatrix} 0 \\ 1 \end{pmatrix}^T (z - \mu) - \frac{1}{2} \begin{pmatrix} 0 \\ 1 \end{pmatrix}^T a \\ &= z_2 - \mu_2 - \frac{1}{2}\sigma_2^2 = z_2. \end{aligned}$$

□

Sketch of proof of Le Cam's 3rd Lemma. Set

$$\Lambda_n := \sum_{i=1}^n \left[\log p_{\theta_n}(X_i) - \log p_{\theta}(X_i) \right].$$

Then under \mathbb{P}_{θ} , by a two-term Taylor expansion,

$$\begin{aligned} \Lambda_n &\approx \frac{h}{\sqrt{n}} \sum_{i=1}^n s_{\theta}(X_i) + \frac{h^2}{2} \frac{1}{n} \sum_{i=1}^n \dot{s}_{\theta}(X_i) \\ &\approx \frac{h}{\sqrt{n}} \sum_{i=1}^n s_{\theta}(X_i) - \frac{h^2}{2} I(\theta), \end{aligned}$$

as

$$\frac{1}{n} \sum_{i=1}^n \dot{s}_{\theta}(X_i) \approx E_{\theta} \dot{s}_{\theta}(X) = -I(\theta).$$

We moreover have, by the assumed asymptotic linearity, under \mathbb{P}_{θ} ,

$$\sqrt{n}(T_n - \theta) \approx \frac{1}{\sqrt{n}} \sum_{i=1}^n l_{\theta}(X_i).$$

Thus,

$$\begin{pmatrix} \sqrt{n}(T_n - \theta) \\ \Lambda_n \end{pmatrix} \xrightarrow{D_{\theta}} Z,$$

where $Z \in \mathbb{R}^2$, has the two-dimensional normal distribution:

$$Z = \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ -\frac{h^2}{2} I(\theta) \end{pmatrix}, \begin{pmatrix} V_{\theta} & hP_{\theta}(l_{\theta}s_{\theta}) \\ hP_{\theta}(l_{\theta}s_{\theta}) & h^2 I(\theta) \end{pmatrix} \right).$$

Thus, we know that for all bounded and continuous $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, one has

$$\mathbf{E}_{\theta} f(\sqrt{n}(T_n - \theta), \Lambda_n) \rightarrow \mathbf{E} f(Z_1, Z_2).$$

Now, let $f : \mathbb{R} \rightarrow \mathbb{R}$ be bounded and continuous. Then, since

$$\prod_{i=1}^n p_{\theta_n}(X_i) = \prod_{i=1}^n p_{\theta}(X_i) e^{\Lambda_n},$$

we may write

$$\mathbf{E}_{\theta_n} f(\sqrt{n}(T_n - \theta)) = \mathbf{E}_{\theta} f(\sqrt{n}(T_n - \theta)) e^{\Lambda_n}.$$

The function $(z_1, z_2) \mapsto f(z_1) e^{z_2}$ is continuous, but not bounded. However, one can show that one may extend the Portmanteau Theorem to this situation. This then yields

$$\mathbf{E}_{\theta} f(\sqrt{n}(T_n - \theta)) e^{\Lambda_n} \rightarrow \mathbf{E} f(Z_1) e^{Z_2}.$$

Now, apply the auxiliary Lemma, with

$$\mu = \begin{pmatrix} 0 \\ -\frac{h^2}{2} I(\theta) \end{pmatrix}, \quad \Sigma = \begin{pmatrix} V_{\theta} & hP_{\theta}(l_{\theta}s_{\theta}) \\ hP_{\theta}(l_{\theta}s_{\theta}) & h^2 I(\theta) \end{pmatrix}.$$

Then we get

$$\mathbf{E}f(Z_1)e^{Z_2} = \int f(z_1)e^{z_2}\phi_Z(z)dz = \int f(z_1)\phi_Y(z)dz = \mathbf{E}f(Y_1),$$

where

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} hP_\theta(l_\theta s_\theta) \\ \frac{h^2}{2}I(\theta) \end{pmatrix}, \begin{pmatrix} V_\theta & hP_\theta(l_\theta s_\theta) \\ hP_\theta(l_\theta s_\theta) & h^2I(\theta) \end{pmatrix}\right),$$

so that

$$Y_1 \sim \mathcal{N}(hP_\theta(l_\theta s_\theta), V_\theta).$$

So we conclude that

$$\sqrt{n}(T_n - \theta) \xrightarrow{D_{\theta_n}} Y_1 \sim \mathcal{N}(hP_\theta(l_\theta s_\theta), V_\theta).$$

Hence

$$\sqrt{n}(T_n - \theta_n) = \sqrt{n}(T_n - \theta) - h \xrightarrow{D_{\theta_n}} \mathcal{N}(h\{P_\theta(l_\theta s_\theta) - 1\}, V_\theta).$$

□

6.7 Asymptotic confidence intervals and tests

Again throughout this section, enough regularity is assumed, such as existence of derivatives and interchanging integration and differentiation.



Intermezzo: the χ^2 distribution Let Y_1, \dots, Y_p be i.i.d. $\mathcal{N}(0, 1)$ -distributed. Define the p -vector

$$Y := \begin{pmatrix} Y_1 \\ \vdots \\ Y_p \end{pmatrix}.$$

Then Y is $\mathcal{N}(0, I)$ -distributed, with I the $p \times p$ identity matrix. The χ^2 -distribution with p degrees of freedom is defined as the distribution of

$$\|Y\|^2 := \sum_{j=1}^p Y_j^2.$$

Notation: $\|Y\|^2 \sim \chi_p^2$.

For a symmetric positive definite matrix Σ , one can define the square root $\Sigma^{1/2}$ as a symmetric positive definite matrix satisfying

$$\Sigma^{1/2}\Sigma^{1/2} = \Sigma.$$

Its inverse is denoted by $\Sigma^{-1/2}$ (which is the square root of Σ^{-1}). If $Z \in \mathbb{R}^p$ is $\mathcal{N}(0, \Sigma)$ -distributed, the transformed vector

$$Y := \Sigma^{-1/2}Z$$

is $\mathcal{N}(0, I)$ -distributed. It follows that

$$Z^T \Sigma^{-1} Z = Y^T Y = \|Y\|^2 \sim \chi_p^2.$$

Asymptotic pivots Recall the definition of an asymptotic pivot (see Section 1.7). It is a function $Z_n(\gamma) := Z_n(X_1, \dots, X_n, \gamma)$ of the data X_1, \dots, X_n and the parameter of interest $\gamma = g(\theta) \in \mathbb{R}^p$, such that its asymptotic distribution does not on the unknown parameter θ , i.e., for a random variable Z , with distribution Q not depending on θ ,

$$Z_n(\gamma) \xrightarrow{D_\theta} Z, \quad \forall \theta.$$

An asymptotic pivot can be used to construct approximate $(1 - \alpha)$ -confidence intervals for γ , and tests for $H_0 : \gamma = \gamma_0$ with approximate level α .

Consider now an asymptotically normal estimator T_n of γ , which is asymptotically unbiased and has asymptotic covariance matrix V_θ , that is

$$\sqrt{n}(T_n - \gamma) \xrightarrow{D_\theta} \mathcal{N}(0, V_\theta), \quad \forall \theta.$$

(assuming such an estimator exists). Then, depending on the situation, there are various ways to construct an asymptotic pivot.

1st asymptotic pivot

If the asymptotic covariance matrix V_θ is non-singular, and depends only on the parameter of interest γ , say $V_\theta = V(\gamma)$ (for example, if $\gamma = \theta$), then an asymptotic pivot is

$$Z_{n,1}(\gamma) := n(T_n - \gamma)^T V(\gamma)^{-1} (T_n - \gamma).$$

The asymptotic distribution is the χ^2 -distribution with p degrees of freedom.

2nd asymptotic pivot

If, for all θ , one has a consistent estimator \hat{V}_n of $V(\theta)$, then an asymptotic pivot is

$$Z_{n,2}(\gamma) := n(T_n - \gamma)^T \hat{V}_n^{-1} (T_n - \gamma).$$

The asymptotic distribution is again the χ^2 -distribution with p degrees of freedom.

Estimators of the asymptotic variance

◦ If $\hat{\theta}_n$ is a consistent estimator of θ and if $\theta \mapsto V_\theta$ is continuous, one may insert $\hat{V}_n := V_{\hat{\theta}_n}$.

◦ If $T_n = \hat{\gamma}_n$ is the M-estimator of γ , γ being the solution of $P_\theta \psi_\gamma = 0$, then (under regularity) the asymptotic covariance matrix is

$$V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1},$$

where

$$J_\theta = P_\theta \psi_\gamma \psi_\gamma^T,$$

and

$$M_\theta = \left. \frac{\partial}{\partial c^T} P_\theta \psi_c \right|_{c=\gamma} = P_\theta \dot{\psi}_\gamma.$$

Then one may estimate J_θ and M_θ by

$$\hat{J}_n := \hat{P}_n \psi_{\hat{\gamma}_n} \psi_{\hat{\gamma}_n}^T = \frac{1}{n} \sum_{i=1}^n \psi_{\hat{\gamma}_n}(X_i) \psi_{\hat{\gamma}_n}^T(X_i),$$

and

$$\hat{M}_n := \hat{P}_n \dot{\psi}_{\hat{\gamma}_n} = \frac{1}{n} \sum_{i=1}^n \dot{\psi}_{\hat{\gamma}_n}(X_i),$$

respectively. Under some regularity conditions,

$$\hat{V}_n := \hat{M}_n^{-1} \hat{J}_n \hat{M}_n^{-1}.$$

is a consistent estimator of V_θ ⁶.

6.7.1 Maximum likelihood

Suppose now that $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ has $\Theta \subset \mathbb{R}^p$, and that \mathcal{P} is dominated by some σ -finite measure ν . Let $p_\theta := dP_\theta/d\nu$ denote the densities, and let

$$\hat{\theta}_n := \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_\vartheta(X_i)$$

be the MLE. Recall that $\hat{\theta}_n$ is an M-estimator with loss function $\rho_\theta = -\log p_\theta$, and hence (under regularity conditions), $\psi_\theta = \dot{\rho}_\theta$ is minus the score function $s_\theta := \dot{p}_\theta/p_\theta$. The asymptotic variance of the MLE is $I^{-1}(\theta)$, where $I(\theta) := P_\theta s_\theta s_\theta^T$ is the Fisher information:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D_\theta} \mathcal{N}(0, I^{-1}(\theta)), \quad \forall \theta.$$

Thus, in this case

$$Z_{n,1}(\theta) = n(\hat{\theta}_n - \theta)I(\theta)(\hat{\theta}_n - \theta),$$

and, with \hat{I}_n being a consistent estimator of $I(\theta)$

$$Z_{n,2}(\theta) = n(\hat{\theta}_n - \theta)\hat{I}_n(\hat{\theta}_n - \theta).$$

⁶From most algorithms used to compute the M-estimator $\hat{\gamma}_n$, one easily can obtain \hat{M}_n and \hat{J}_n as output. Recall e.g. that the Newton-Raphson algorithm is based on the iterations

$$\hat{\gamma}_{\text{new}} = \hat{\gamma}_{\text{old}} - \left(\sum_{i=1}^n \dot{\psi}_{\hat{\gamma}_{\text{old}}} \right)^{-1} \sum_{i=1}^n \psi_{\hat{\gamma}_{\text{old}}}.$$

Note that one may take

$$\hat{I}_n := -\frac{1}{n} \sum_{i=1}^n \dot{s}_{\hat{\theta}_n}(X_i) = -\frac{\partial^2}{\partial \vartheta \partial \vartheta^T} \frac{1}{n} \sum_{i=1}^n \log p_{\vartheta}(X_i) \Big|_{\vartheta=\hat{\theta}_n}$$

as estimator of the Fisher information ⁷.

3rd asymptotic pivot

Define now the twice log-likelihood ratio

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) := 2 \sum_{i=1}^n \left[\log p_{\hat{\theta}_n}(X_i) - \log p_{\theta}(X_i) \right].$$

It turns out that the log-likelihood ratio is indeed an asymptotic pivot. A practical advantage is that it is self-normalizing: one does not need to explicitly estimate asymptotic (co-)variances.

Lemma 6.7.1 *Under regularity conditions, $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta)$ is an asymptotic pivot for θ . Its asymptotic distribution is again the χ^2 -distribution with p degrees of freedom:*

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \xrightarrow{D_{\theta}} \chi_p^2 \quad \forall \theta.$$

Sketch of the proof. We have by a two-term Taylor expansion

$$\begin{aligned} 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) &= 2n\hat{P}_n \left[\log p_{\hat{\theta}_n} - \log p_{\theta} \right] \\ &\approx 2n(\hat{\theta}_n - \theta)^T \hat{P}_n s_{\theta} + n(\hat{\theta}_n - \theta)^T \hat{P}_n \dot{s}_{\theta}(\hat{\theta}_n - \theta) \\ &\approx 2n(\hat{\theta}_n - \theta)^T \hat{P}_n s_{\theta} - n(\hat{\theta}_n - \theta)^T I(\theta)(\hat{\theta}_n - \theta), \end{aligned}$$

where in the second step, we used $\hat{P}_n \dot{s}_{\theta} \approx P_{\theta} \dot{s}_{\theta} = -I(\theta)$. (You may compare this two-term Taylor expansion with the one in the sketch of proof of Le Cam's 3rd Lemma). The MLE $\hat{\theta}_n$ is asymptotically linear with influence function $l_{\theta} = I(\theta)^{-1} s_{\theta}$:

$$\hat{\theta}_n - \theta = I(\theta)^{-1} \hat{P}_n s_{\theta} + o_{\mathbf{P}_{\theta}}(n^{-1/2}).$$

Hence,

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \approx n(\hat{P}_n s_{\theta})^T I(\theta)^{-1} (\hat{P}_n s_{\theta}).$$

The result now follows from

$$\sqrt{n} \hat{P}_n s_{\theta} \xrightarrow{D_{\theta}} \mathcal{N}(0, I(\theta)).$$

□

⁷In other words (as for general M-estimators), the algorithm (e.g. Newton Raphson) for calculating the maximum likelihood estimator $\hat{\theta}_n$ generally also provides an estimator of the Fisher information as by-product.

Example 6.7.1 Let X_1, \dots, X_n be i.i.d. copies of X , where $X \in \{1, \dots, k\}$ is a label, with

$$P_\theta(X = j) := \pi_j, \quad j = 1, \dots, k.$$

where the probabilities π_j are positive and add up to one: $\sum_{j=1}^k \pi_j = 1$, but are assumed to be otherwise unknown. Then there are $p := k - 1$ unknown parameters, say $\theta = (\pi_1, \dots, \pi_{k-1})$. Define $N_j := \#\{i : X_i = j\}$. (Note that (N_1, \dots, N_k) has a multinomial distribution with parameters n and (π_1, \dots, π_k)).

Lemma For each $j = 1, \dots, k$, the MLE of π_j is

$$\hat{\pi}_j = \frac{N_j}{n}.$$

Proof. The log-densities can be written as

$$\log p_\theta(x) = \sum_{j=1}^k 1\{x = j\} \log \pi_j,$$

so that

$$\sum_{i=1}^n \log p_\theta(X_i) = \sum_{j=1}^k N_j \log \pi_j.$$

Putting the derivatives with respect to $\theta = (\pi_1, \dots, \pi_{k-1})$, (with $\pi_k = 1 - \sum_{j=1}^{k-1} \pi_j$) to zero gives,

$$\frac{N_j}{\hat{\pi}_j} - \frac{N_k}{\hat{\pi}_k} = 0.$$

Hence

$$\hat{\pi}_j = N_j \frac{\hat{\pi}_k}{N_k}, \quad j = 1, \dots, k,$$

and thus

$$1 = \sum_{j=1}^k \hat{\pi}_j = n \frac{\hat{\pi}_k}{N_k},$$

yielding

$$\hat{\pi}_k = \frac{N_k}{n},$$

and hence

$$\hat{\pi}_j = \frac{N_j}{n}, \quad j = 1, \dots, k.$$

□

We now first calculate $Z_{n,1}(\theta)$. For that, we need to find the Fisher information $I(\theta)$.

Lemma *The Fisher information is*

$$I(\theta) = \begin{pmatrix} \frac{1}{\pi_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\pi_{k-1}} \end{pmatrix} + \frac{1}{\pi_k} \iota \iota^T, \quad ^8$$

where ι is the $(k-1)$ -vector $\iota := (1, \dots, 1)^T$.

Proof. We have

$$s_{\theta,j}(x) = \frac{1}{\pi_j} \mathbf{1}\{x = j\} - \frac{1}{\pi_k} \mathbf{1}\{x = k\}.$$

So

$$\begin{aligned} (I(\theta))_{j_1, j_2} &= E_\theta \left(\frac{1}{\pi_{j_1}} \mathbf{1}\{X = j_1\} - \frac{1}{\pi_k} \mathbf{1}\{X = k\} \right) \left(\frac{1}{\pi_{j_2}} \mathbf{1}\{X = j_2\} - \frac{1}{\pi_k} \mathbf{1}\{X = k\} \right) \\ &= \begin{cases} \frac{1}{\pi_k} & j_1 \neq j_2 \\ \frac{1}{\pi_j} + \frac{1}{\pi_k} & j_1 = j_2 = j \end{cases}. \end{aligned}$$

□

We thus find

$$\begin{aligned} Z_{n,1}(\theta) &= n(\hat{\theta}_n - \theta)^T I(\theta) (\hat{\theta}_n - \theta) \\ &= n \begin{pmatrix} \hat{\pi}_1 - \pi_1 \\ \vdots \\ \hat{\pi}_{k-1} - \pi_{k-1} \end{pmatrix}^T \left(\begin{pmatrix} \frac{1}{\pi_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\pi_{k-1}} \end{pmatrix} + \frac{1}{\pi_k} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} \right) \begin{pmatrix} \hat{\pi}_1 - \pi_1 \\ \vdots \\ \hat{\pi}_{k-1} - \pi_{k-1} \end{pmatrix} \\ &= n \sum_{j=1}^{k-1} \frac{(\hat{\pi}_j - \pi_j)^2}{\pi_j} + n \frac{1}{\pi_k} \left(\sum_{j=1}^{k-1} (\hat{\pi}_j - \pi_j) \right)^2 \\ &= n \sum_{j=1}^k \frac{(\hat{\pi}_j - \pi_j)^2}{\pi_j} \\ &= \sum_{j=1}^k \frac{(N_j - n\pi_j)^2}{n\pi_j}. \end{aligned}$$

This is called the Pearson's chi-square

$$\sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}.$$

A version of $Z_{n,2}(\theta)$ is to replace, for $j = 1, \dots, k$, π_j by $\hat{\pi}_j$ in the expression for the Fisher information. This gives

$$Z_{n,2}(\theta) = \sum_{j=1}^k \frac{(N_j - n\pi_j)^2}{N_j}.$$

⁸To invert such a matrix, one may apply the formula $(A + bb^T)^{-1} = A^{-1} - \frac{A^{-1}bb^T A^{-1}}{1 + b^T A^{-1}b}$.

This is called the Pearson's chi-square

$$\sum \frac{(\text{observed} - \text{expected})^2}{\text{observed}}.$$

Finally, the log-likelihood ratio pivot is

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) = 2 \sum_{j=1}^k N_j \log \left(\frac{\hat{\pi}_j}{\pi_j} \right).$$

The approximation $\log(1+x) \approx x - x^2/2$ shows that $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \approx Z_{n,2}(\theta)$:

$$\begin{aligned} 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) &= -2 \sum_{j=1}^k N_j \log \left(1 + \frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right) \\ &\approx -2 \sum_{j=1}^k N_j \left(\frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right) + \sum_{j=1}^k N_j \left(\frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right)^2 = Z_{n,2}(\theta). \end{aligned}$$

The three asymptotic pivots $Z_{n,1}(\theta)$, $Z_{n,2}(\theta)$ and $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta)$ are each asymptotically χ_{k-1}^2 -distributed under \mathbb{P}_θ .

6.7.2 Likelihood ratio tests

Intermezzo: some matrix algebra

Let $z \in \mathbb{R}^p$ be a vector and B be a $(q \times p)$ -matrix, $(p \geq q)$ with rank q . Moreover, let V be a positive definite $(p \times p)$ -matrix.

Lemma *We have*

$$\max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T a\} = z^T z - z^T B^T (BB^T)^{-1} Bz.$$

Proof. We use Lagrange multipliers $\lambda \in \mathbb{R}^p$. We have

$$\frac{\partial}{\partial a} \{2a^T z - a^T a + 2a^T B^T \lambda\} = z - a + B^T \lambda.$$

Hence for

$$a_* := \arg \max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T a\},$$

we have

$$z - a_* + B^T \lambda = 0,$$

or

$$a_* = z + B^T \lambda.$$

The restriction $Ba_* = 0$ gives

$$Bz + BB^T \lambda = 0.$$

So

$$\lambda = -(BB^T)^{-1}Bz.$$

Inserting this in the solution a^* gives

$$a_* = z - B^T(BB^T)^{-1}Bz.$$

Now,

$$a_*^T a_* = (z^T - z^T B^T (BB^T)^{-1} B)(z - B^T (BB^T)^{-1} Bz) = z^T z - z^T B^T (BB^T)^{-1} Bz.$$

So

$$2a_*^T z - a_*^T a_* = z^T z - z^T B^T (BB^T)^{-1} Bz.$$

□

Lemma *We have*

$$\max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T Va\} = z^T V^{-1} z - z^T V^{-1} B^T (BV^{-1} B^T)^{-1} BV^{-1} z.$$

Proof. Make the transformation $b := V^{1/2}a$, and $y := V^{-1/2}z$, and $C = BV^{-1/2}$. Then

$$\begin{aligned} & \max_{a: Ba=0} \{2a^T z - a^T Va\} \\ &= \max_{b: Cb=0} \{2b^T y - b^T b\} \\ &= y^T y - y^T C^T (CC^T)^{-1} C y = z^T V^{-1} z - z^T V^{-1} B^T (BV^{-1} B^T)^{-1} BV^{-1} z. \end{aligned}$$

□

Corollary *Let $L(a) := 2a^T z - a^T Va$. The difference between the unrestricted maximum and the restricted maximum of $L(a)$ is*

$$\max_a L(a) - \max_{a: Ba=0} L(a) = z^T V^{-1} B^T (BV^{-1} B^T)^{-1} BV^{-1} z.$$

Hypothesis testing

For the simple hypothesis

$$H_0 : \theta = \theta_0,$$

we can use $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta_0)$ as test statistic: reject H_0 if $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta_0) > \chi_{p,\alpha}^2$, where $\chi_{p,\alpha}^2$ is the $(1 - \alpha)$ -quantile of the χ_p^2 -distribution.

Consider now the hypothesis

$$H_0 : R(\theta) = 0,$$

where

$$R(\theta) = \begin{pmatrix} R_1(\theta) \\ \vdots \\ R_q(\theta) \end{pmatrix}.$$

Let $\hat{\theta}_n$ be the unrestricted MLE, that is

$$\hat{\theta}_n = \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_{\vartheta}(X_i).$$

Moreover, let $\hat{\theta}_n^0$ be the restricted MLE, defined as

$$\hat{\theta}_n^0 = \arg \max_{\vartheta \in \Theta: R(\vartheta)=0} \sum_{i=1}^n \log p_{\vartheta}(X_i).$$

Define the $(q \times p)$ -matrix

$$\dot{R}(\theta) = \frac{\partial}{\partial \vartheta^T} R(\vartheta)|_{\vartheta=\theta}.$$

We assume $\dot{R}(\theta)$ has rank q .

Let

$$\mathcal{L}_n(\hat{\theta}_n) - \mathcal{L}_n(\hat{\theta}_n^0) = \sum_{i=1}^n \left[\log p_{\hat{\theta}_n}(X_i) - \log p_{\hat{\theta}_n^0}(X_i) \right]$$

be the log-likelihood ratio for testing $H_0 : R(\theta) = 0$.

Lemma 6.7.2 *Under regularity conditions, and if $H_0 : R(\theta) = 0$ holds, we have*

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \xrightarrow{D_{\theta}} \chi_q^2.$$

Sketch of the proof. Let

$$\mathbf{Z}_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n s_{\theta}(X_i).$$

As in the sketch of the proof of Lemma 6.7.1, we can use a two-term Taylor expansion to show for any sequence ϑ_n satisfying $\vartheta_n = \theta + O_{\mathbf{P}_{\theta}}(n^{-1/2})$, that

$$2 \sum_{i=1}^n \left[\log p_{\vartheta_n}(X_i) - \log p_{\theta}(X_i) \right] = 2\sqrt{n}(\vartheta_n - \theta)^T \mathbf{Z}_n - n(\vartheta_n - \theta)^2 I(\theta)(\vartheta_n - \theta) + o_{\mathbf{P}_{\theta}}(1).$$

Here, we also again use that $\sum_{i=1}^n \dot{s}_{\vartheta_n}(X_i)/n = -I(\theta) + o_{\mathbf{P}_{\theta}}(1)$. Moreover, by a one-term Taylor expansion, and invoking that $R(\theta) = 0$,

$$R(\vartheta_n) = \dot{R}(\theta)(\vartheta_n - \theta) + o_{\mathbf{P}_{\theta}}(n^{-1/2}).$$

Insert the corollary in the above matrix algebra, with $z := \mathbf{Z}_n$, $B := \dot{R}(\theta)$, and $V = I(\theta)$. This gives

$$\begin{aligned} & 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \\ &= 2 \sum_{i=1}^n \left[\log p_{\hat{\theta}_n}(X_i) - \log p_{\theta}(X_i) \right] - 2 \sum_{i=1}^n \left[\log p_{\hat{\theta}_n^0}(X_i) - \log p_{\theta}(X_i) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbf{Z}_n^T I(\theta)^{-1} \dot{R}^T(\theta) \left(\dot{R}(\theta) I(\theta)^{-1} \dot{R}(\theta)^T \right)^{-1} \dot{R}(\theta) I(\theta)^{-1} \mathbf{Z}_n + o_{\mathbf{P}_\theta}(1) \\
&:= \mathbf{Y}_n^T W^{-1} \mathbf{Y}_n + o_{\mathbf{P}_\theta}(1),
\end{aligned}$$

where \mathbf{Y}_n is the q -vector

$$\mathbf{Y}_n := \dot{R}(\theta) I(\theta)^{-1} \mathbf{Z}_n,$$

and where W is the $(q \times q)$ -matrix

$$W := \dot{R}(\theta) I(\theta)^{-1} \dot{R}(\theta)^T.$$

We know that

$$\mathbf{Z}_n \xrightarrow{D_\theta} \mathcal{N}(0, I(\theta)).$$

Hence

$$\mathbf{Y}_n \xrightarrow{D_\theta} \mathcal{N}(0, W),$$

so that

$$\mathbf{Y}_n^T W^{-1} \mathbf{Y}_n \xrightarrow{D_\theta} \chi_q^2.$$

□

Corollary 6.7.1 *From the sketch of the proof of Lemma 6.7.2, one sees that moreover (under regularity),*

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n(\hat{\theta}_n - \hat{\theta}_n^0)^T I(\theta)(\hat{\theta}_n - \hat{\theta}_n^0),$$

and also

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n(\hat{\theta}_n - \hat{\theta}_n^0)^T I(\hat{\theta}_n^0)(\hat{\theta}_n - \hat{\theta}_n^0).$$

Example 6.7.2 Let X be a bivariate label, say $X \in \{(j, k) : j = 1, \dots, r, k = 1, \dots, s\}$. For example, the first index may correspond to sex ($r = 2$) and the second index to the color of the eyes ($s = 3$). The probability of the combination (j, k) is

$$\pi_{j,k} := P_\theta \left(X = (j, k) \right).$$

Let X_1, \dots, X_n be i.i.d. copies of X , and

$$N_{j,k} := \#\{X_i = (j, k)\}.$$

From Example 6.7.1, we know that the (unrestricted) MLE of $\pi_{j,k}$ is equal to

$$\hat{\pi}_{j,k} := \frac{N_{j,k}}{n}.$$

We now want to test whether the two labels are independent. The null-hypothesis is

$$H_0 : \pi_{j,k} = (\pi_{j,+}) \times (\pi_{+,k}) \quad \forall (j, k).$$

Here

$$\pi_{j,+} := \sum_{k=1}^s \pi_{j,k}, \quad \pi_{+,k} := \sum_{j=1}^r \pi_{j,k}.$$

One may check that the restricted MLE is

$$\hat{\pi}_{j,k}^0 = (\hat{\pi}_{j,+}) \times (\hat{\pi}_{+,k}),$$

where

$$\hat{\pi}_{j,+} := \sum_{k=1}^s \hat{\pi}_{j,k}, \quad \hat{\pi}_{+,k} := \sum_{j=1}^r \hat{\pi}_{j,k}.$$

The log-likelihood ratio test statistic is thus

$$\begin{aligned} 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) &= 2 \sum_{j=1}^r \sum_{k=1}^s N_{j,k} \left[\log \left(\frac{N_{j,k}}{n} \right) - \log \left(\frac{N_{j,+} N_{+,k}}{n^2} \right) \right] \\ &= 2 \sum_{j=1}^r \sum_{k=1}^s N_{j,k} \log \left(\frac{n N_{j,k}}{N_{j,+} N_{+,k}} \right). \end{aligned}$$

Its approximation as given in Corollary 6.7.1 is

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n \sum_{j=1}^r \sum_{k=1}^s \frac{(N_{j,k} - N_{j,+} N_{+,k}/n)^2}{N_{j,+} N_{+,k}}.$$

This is Pearson's chi-squared test statistic for testing independence. To find out what the value of q is in this example, we first observe that the unrestricted case has $p = rs - 1$ free parameters. Under the null-hypothesis, there remain $(r - 1) + (s - 1)$ free parameters. Hence, the number of restrictions is

$$q = (rs - 1) - ((r - 1) + (s - 1)) = (r - 1)(s - 1).$$

Thus, under $H_0 : \pi_{j,k} = (\pi_{j,+}) \times (\pi_{+,k}) \forall (j, k)$, we have

$$n \sum_{j=1}^r \sum_{k=1}^s \frac{(N_{j,k} - N_{j,+} N_{+,k}/n)^2}{N_{j,+} N_{+,k}} \xrightarrow{D_\theta} \chi_{(r-1)(s-1)}^2.$$

6.8 Complexity regularization (to be written)

Chapter 7

Literature

- J.O. Berger (1985) *Statistical Decision Theory and Bayesian Analysis* Springer
A fundamental book on Bayesian theory.
- P.J. Bickel, K.A. Doksum (2001) *Mathematical Statistics, Basic Ideas and Selected Topics* Volume I, 2nd edition, Prentice Hall
Quite general, and mathematically sound.
- D.R. Cox and D.V. Hinkley (1974) *Theoretical Statistics* Chapman and Hall
Contains good discussions of various concepts and their practical meaning. Mathematical development is sketchy.
- J.G. Kalbfleisch (1985) *Probability and Statistical Inference* Volume 2, Springer
Treats likelihood methods.
- L.M. Le Cam (1986) *Asymptotic Methods in Statistical Decision Theory* Springer
Treats decision theory on a very abstract level.
- E.L. Lehmann (1983) *Theory of Point Estimation* Wiley
A “klassiker”. The lecture notes partly follow this book
- E.L. Lehmann (1986) *Testing Statistical Hypothesis* 2nd edition, Wiley
Goes with the previous book.
- J.A. Rice (1994) *Mathematical Statistics and Data Analysis* 2nd edition, Duxbury Press
A more elementary book.
- M.J. Schervish (1995) *Theory of Statistics* Springer
Mathematically exact and quite general. Also good as reference book.
- R.J. Serfling (1980) *Approximation Theorems of Mathematical Statistics* Wiley

Treats asymptotics.

- A.W. van der Vaart (1998) *Asymptotic Statistics* Cambridge University Press
Treats modern asymptotics and e.g. semiparametric theory