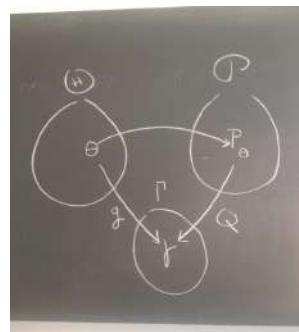


Mathematical Statistics



Sara van de Geer

September 2018

Contents

1	Introduction	9
1.1	Speed of light example	9
1.2	Notation	10
1.3	Example: the location model	11
1.4	Some further examples of statistical models	12
2	Estimation	15
2.1	What is an estimator?	15
2.2	The empirical distribution function	15
2.3	Some estimators for the location model	16
2.4	How to construct estimators	17
2.4.1	Plug-in estimators	17
2.4.2	The method of moments	18
2.4.3	Likelihood methods	20
2.5	Asymptotic tests and confidence intervals based on the likelihood	23
3	Intermezzo: distribution theory	25
3.1	Conditional distributions	25
3.2	The multinomial distribution	26
3.3	The Poisson distribution	27
3.4	The distribution of the maximum of two random variables	28
4	Sufficiency and exponential families	31
4.1	Sufficiency	31
4.2	Factorization Theorem of Neyman	33
4.3	Exponential families	34
4.4	Intermezzo: the mean and covariance matrix of a random vector	36
4.5	Canonical form of an exponential family	37
4.6	Reparametrizing in the one-dimensional case	38
4.7	Score function and Fisher information	39
4.8	Score function for exponential families	40
4.9	Minimal sufficiency	42
5	Bias, variance and the Cramér Rao lower bound	43
5.1	What is an unbiased estimator?	43
5.2	UMVU estimators	44

5.3	The Lehmann-Scheff� Lemma	47
5.4	Completeness for exponential families	49
5.5	The Cram� Rao lower bound	51
5.6	CRLB and exponential families	54
5.7	Higher-dimensional extensions	55
6	Tests and confidence intervals	59
6.1	Intermezzo: quantile functions	59
6.2	How to construct tests	59
6.3	Equivalence confidence sets and tests	61
6.4	Comparison of confidence intervals and tests	62
6.5	An illustration: the two-sample problem	62
6.5.1	Student's test	63
6.5.2	Wilcoxon's test	64
6.5.3	Comparison of Student's test and Wilcoxon's test	67
7	The Neyman Pearson Lemma and UMP tests	69
7.1	The Neyman Pearson Lemma	69
7.2	Uniformly most powerful tests	70
7.2.1	An example	70
7.3	UMP tests and exponential families	72
7.4	One- and two-sided tests: an example with the Bernoulli distribution	74
7.5	Unbiased tests	75
7.6	Conditional tests \star	78
8	Comparison of estimators	83
8.1	Definition of risk	83
8.2	Risk and sufficiency	84
8.3	Rao-Blackwell	84
8.4	Sensitivity and robustness	85
8.5	Computational aspects	86
9	Equivariant statistics	87
9.1	Equivariance in the location model	87
9.1.1	Construction of the UMRE estimator	88
9.1.2	Quadratic loss: the Pitman estimator	89
9.1.3	Invariant statistics	91
9.1.4	Quadratic loss and Basu's Lemma	91
9.2	Equivariance in the location-scale model \star	93
9.2.1	Construction of the UMRE estimator \star	94
9.2.2	Quadratic loss \star	94
10	Decision theory	97
10.1	Decisions and their risk	97
10.2	Admissible decisions	99
10.2.1	Not using the data at all is admissible	99

10.2.2 A Neyman Pearson test is admissible	100
10.3 Minimax decisions	101
10.3.1 Minimax Neyman Pearson test	101
10.4 Bayes decisions	101
10.4.1 Bayes test	102
10.5 Construction of Bayes estimators	103
10.5.1 Bayes test revisited	104
10.5.2 Bayes estimator for quadratic loss	104
10.5.3 Bayes estimator and the maximum a posteriori estimator	105
10.5.4 Three worked-out examples	105
10.6 Discussion of Bayesian approach	107
10.7 Integrating parameters out \star	109
11 Proving admissibility and minimaxity	111
11.1 Minimaxity	112
11.1.1 Minimaxity of the Pitman estimator \star	113
11.2 Admissibility	114
11.2.1 Admissible estimators for the normal mean	116
11.3 Admissible estimators in exponential families \star	118
11.4 Inadmissibility in higher-dimensional settings \star	120
12 The linear model	123
12.1 Definition of the least squares estimator	123
12.2 Intermezzo: the χ^2 distribution	126
12.3 Distribution of the least squares estimator	127
12.4 Intermezzo: some matrix algebra	128
12.5 Testing a linear hypothesis	130
13 Asymptotic theory	131
13.1 Types of convergence	132
13.1.1 Stochastic order symbols	133
13.1.2 Some implications of convergence	134
13.2 Consistency and asymptotic normality	135
13.3 Asymptotic linearity	135
13.4 The δ -technique	137
14 M-estimators	141
14.1 MLE as special case of M-estimation	142
14.2 Consistency of M-estimators	144
14.3 Asymptotic normality of M-estimators	147
14.4 Asymptotic normality of the MLE	150
14.5 Two further examples of M-estimation	151
14.6 Asymptotic relative efficiency	154
14.7 Asymptotic pivots	156
14.8 Asymptotic pivot based on the MLE	157
14.9 MLE for the multinomial distribution	159
14.10 Likelihood ratio tests	161

14.11 Contingency tables	164
15 Abstract asymptotics *	167
15.1 Plug-in estimators *	168
15.2 Consistency of plug-in estimators *	170
15.3 Asymptotic normality of plug-in estimators *	171
15.4 Asymptotic Cramer Rao lower bound *	175
15.5 Le Cam's 3 rd Lemma *	177
16 Complexity regularization	183
16.1 Non-parametric regression	183
16.2 Smoothness classes	184
16.3 A continuous version with explicit solution *	185
16.4 Estimating a function of bounded variation	186
16.5 The ridge and Lasso penalty	188
16.6 Conclusion	193
17 Literature	195

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

Throughout the notes measurability assumptions are not stated explicitly.

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

Chapters or (sub)sections with a \star are not part of the exam.



Chapter 1

Introduction

Statistics is about the mathematical modeling of observable phenomena, using stochastic models, and about analyzing data: estimating parameters of the model, constructing confidence intervals 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.

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.

1.1 Speed of light example

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.

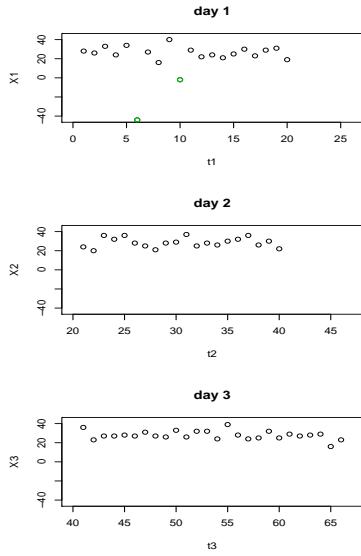
distance: 7.44373 km.

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:



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.

1.2 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. Typically the distributions in \mathcal{P} are indexed by a parameter, say θ , in some parameter space, say Θ . Then $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$, and $P = P_\theta$ for some $\theta \in \Theta$. Then θ is often called the “true parameter”.¹ The parameter space Θ may be high-dimensional, or even ∞ -dimensional. Often, one is only interested in a certain aspect of the parameter. We write the parameter of interest as $\gamma := g(\theta)$ where $g : \Theta \rightarrow \Gamma$ a given function with values in some space Γ .

1.3 Example: the location model

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

Let X be real-valued. The location model is

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

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. Then $g(\mu, F_0) = \mu$.

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 ∞ -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

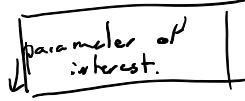
$$\mathcal{F}_0 := \{\Phi(\cdot/\sigma) : \sigma > 0\}, \quad (1.3)$$

where Φ is the standard normal distribution function.

¹To be mathematically correct one should write $P_\theta \in \{P_\vartheta : \vartheta \in \Theta\}$ to make the distinction in notation between the true parameter θ and the parameter ϑ indexing the class \mathcal{P} . We actually need this distinction as the theory develops.

this is
an example
of the nuance
parameters discussed
in Bayesian statistics.

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 .

1.4 Some further examples of statistical models

Example 1.4.1 Poisson distribution

Consider a small insurance company. Let X be the number of claims on a particular day. Then a possible model is the Poisson model, which assumes that X has a Poisson distribution with parameter $\theta > 0$:

$$P_\theta(X = x) = \frac{\theta^x}{x!} e^{-\theta}, \quad x = 0, 1, 2, \dots$$

A parameter of interest could be for example the probability of at least 4 claims on a particular day. Then

$$\begin{aligned} \gamma &= P_\theta(X \geq 4) \\ &= 1 - P_\theta(X \leq 3) \\ &= 1 - \left(1 + \theta + \frac{\theta^2}{2} + \frac{\theta^3}{3!}\right) e^{-\theta} \\ &:= g(\theta). \end{aligned}$$

Suppose we observed the number of claims X_1, \dots, X_n during $n = 200$ days. A possible estimator of θ is the sample average $\bar{X} := \sum_{i=1}^n X_i/n$. For γ we may use the “plug in” principle, that is, we plug in the estimator \bar{X} for θ in the function g . This gives $\hat{\gamma} := g(\bar{X})$ as estimator of $g(\theta)$.

Here is a data example

x_i	# days
0	100
1	60
2	32
3	8
≥ 4	0

The observed average is $\bar{x} = .74$ and the estimate of $P_\theta(X \geq 4)$ is .00697.

Example 1.4.2 Pareto distribution

Suppose X has density (with respect to Lebesgue measure ν)

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

where $\theta > 0$ is unknown. This is the Pareto density which is often used to model the distribution of income. The parameter θ is sometimes called the Pareto index. We write the model class for the distribution as

$$\mathcal{P} = \{P_\theta : dP_\theta/d\nu = p_\theta\}.$$

A parameter of interest may be the Gini index. It describes income inequality and is defined for $\theta \geq 1$ as $\gamma(\theta) = 1/(2\theta - 1)$. The value $G(1) = 1$ for $\theta = 1$ corresponds to complete income inequality.

Example 1.4.3 Classification

Let $X = (Y, Z)$ where $Z = \text{body mass index} \in \mathbb{R}$ and $Y \in \{0, 1\}$ indicates having diabetes or not. We assume the model

$$P_\theta(Y = 1|Z = z) = \theta(z), \quad z \in \mathbb{R},$$

with

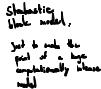
$$\theta(\cdot) \in \Theta := \left\{ \text{all increasing functions } \theta : \mathbb{R} \rightarrow [0, 1] \right\}.$$

The parameter space is ∞ -dimensional. A parameter of interest is for example

$$\gamma := \theta^{-1}\left(\frac{1}{2}\right),$$

that is, the value γ such that for $z \geq \gamma$ you are at risk:

$$P_\theta(Y = 1|Z = z) \geq \frac{1}{2}.$$



⇒ Example 1.4.4 Social networks

Consider p individuals, which either do or do not communicate with each other. If they do, we say there is a connection between them, or we call them friends. Let X be a $p \times p$ matrix $X = (X_{j,k})$ coding the connections: for $j \neq k$

$$X_{j,k} := \begin{cases} 1 & \text{if there is a connection between } j \text{ and } k \\ 0 & \text{else} \end{cases}.$$

The “stochastic block model” assumes that the $\{X_{j,k}\}$ are independent and

$$P_\theta(X_{j,k} = 1) = \begin{cases} \beta_m & \text{if } j \text{ and } k \text{ are in the same community } m \\ \delta & \text{if } j \text{ and } k \text{ are in different communities} \end{cases}.$$

different modeling
depending on whether
in same network or
not. ↗ Network
modeled by
sparse matrix.

If the number of communities is known, say M , but otherwise nothing further, then the parameter space is

$$\Theta = \left\{ (\beta_1, \dots, \beta_M, \delta) \in [0, 1]^{M+1}, \mathcal{M} \right\}$$

\mathcal{M} individuals

where \mathcal{M} is the collection of all community configurations. There are M^p possible community configurations, so $|\mathcal{M}| = M^p$. If p is large this is a huge number, i.e. the parameter space is very complex.

Remark: typically we observe only one realization of X , i.e. $n = 1$.

*If p is large
again because of
the permutations,
this problem becomes
difficult to solve.*

Example 1.4.5 Causal models

Suppose $X = (Z_1, \dots, Z_p) \in \mathbb{R}^p$ is a p -dimensional random variable, for example Z_1 = rainfall, Z_2 = tea consumption, Z_3 = number of tall people, Z_4 = mountain height, ... within a particular canton of Switzerland. A causal model aims to find out which variables are causes and which are consequences.

The structural relations model is

$$Z_{\pi(j)} = f_j \left(Z_{\pi(1)}, \dots, Z_{\pi(j-1)}, \epsilon_j \right), \quad j = 2, \dots, p.$$

Here π is a permutation of $\{1, \dots, p\}$, ϵ_j is unobservable noise and f_j is a partly unknown function. If we assume (for simplicity) the noise distribution to be known the parameter space is

$$\Theta := \left\{ \text{all permutations } \pi \text{ and structural relations } (f_2, \dots, f_p) \right\}.$$

A parameter of interest is for example the causal graph. \Rightarrow 

A sub-example is where the structural relations are modeled as being linear:

$$f_j(z_1, \dots, z_{j-1}, \epsilon_j) = \beta_{j,1}z_1 + \dots + \beta_{j,j-1}z_{j-1} + \epsilon_j, \quad j = 2, \dots, p$$

with $\{\beta_{j,k}\}$ unknown coefficients.

We use the "dt" notation e.g.:

$$F(x) = P(X \leq x) \quad x \in \mathbb{R} \quad \stackrel{\text{short cut}}{\Rightarrow} \quad F = F(\cdot) = P(X \leq \cdot)$$

$$X \sim P := X \text{ has distribution } P.$$

Chapter 2

Estimation

2.1 What is an estimator?

Recall that the data consist of observations $\mathbf{X} = (X_1, \dots, X_n)$ with partly unknown distribution. A parameter is an aspect of the unknown distribution. We typically assume that X_1, \dots, X_n are i.i.d. copies of a random variable X where X has distribution $P_\theta \in \{P_\vartheta : \vartheta \in \Theta\}$.

An estimator is constructed to estimate some unknown parameter, γ say. Its formal definition is

Definition 2.1.1 *An estimator $T(\mathbf{X})$ is some given (measurable) function $T(\cdot)$ evaluated at the observations \mathbf{X} . The function $T(\cdot)$ is not allowed to depend on unknown parameters.*

An estimator is also called a *statistic* or a *decision*.

The reason why T is not allowed to depend on unknown parameters is that one should be able to calculate it in practice, using only the data. We will often use the same notation T for the estimator $T(\mathbf{X})$ (i.e., we write $T = T(\mathbf{X})$ omitting the argument \mathbf{X}) and the function $T = T(\cdot)$. It should be clear from the context which of the two is meant.

2.2 The empirical distribution function

Let X_1, \dots, X_n be real-valued observations. 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

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

Most estimators are constructed according the so-called a **plug-in principle** (*Einsatzprinzip*). 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).

$T = \text{estimator of } \gamma; \text{ of quality}$
of interest

2.3 Some estimators for the location model

In the location model of Section 1.3, one may consider the following estimators $\hat{\mu}$ of μ (among others):

- The average

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

Note that $\hat{\mu}_1$ minimizes over μ the squared loss¹

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

that is

$$\hat{\mu} = \arg \min_{\mu \in \mathbb{R}} \sum_{i=1}^n (X_i - \mu)^2. \quad (2.1)$$

It can be shown that $\hat{\mu}_1$ is a “good” estimator if the model (1.3) holds, i.e., if X_1, \dots, X_n are i.i.d. $\mathcal{N}(\mu, \sigma^2)$. 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

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

¹To avoid misunderstanding, we note that e.g. in (2.1), μ is used as variable over which is minimized and is there not the unknown parameter μ . It is a general convention to abuse notation and employ the same symbol μ . When further developing the theory we shall often introduce a different symbol for the variable, to distinguish it from the “true parameter” μ , e.g., (2.1) is written as

$$\hat{\mu}_1 := \arg \min_{c \in \mathbb{R}} \sum_{i=1}^n (X_i - c)^2.$$

i.e. interpret
of parameter as a value
of the probability
function

So here
is the general
idea of the
chapter.

You have

parameter of
interest

$\gamma = Q(F)$

i.e. depends
on true
function

$F = P_\theta$ and
is a ‘function
of’ a map $Q(\cdot)$
of it.

You then
have an
estimator $\hat{T}(\cdot)$
computed as
 $Q(\hat{F}_n)$ through
plug in!

- The Huber estimator is

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

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.

higher penalty if above threshold.
complex non-linear relation!

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

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

The above estimators $\hat{\mu}_1, \dots, \hat{\mu}_4$ are **plug-in estimators** 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 (i.e. if F is symmetric around μ) $\hat{\mu}_1, \dots, \hat{\mu}_4$ 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 .

2.4 How to construct estimators

2.4.1 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_θ . Hence, a parameter is a function of F_θ :

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

! 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). 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 $Q_n(\hat{P}_n)$.

2.4.2 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

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).

There was an example in St. Gallen.
One had non-invertible moments function.

to exist), i.e.,

defining map $\Omega(F_\theta)$

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

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

plug-in
 $\Omega(F_n)$

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

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

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

*Then invert
if and e
gd of param
of inf inst*

Example 2.4.1 Method of moments for the negative binomial distribution

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

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

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

$$E_\theta(X) = k \frac{(1-\theta)}{\theta} := m(\theta). \quad \begin{aligned} \mu &= k \frac{(1-\theta)}{\theta} \\ \mu\theta &= k - k\theta \end{aligned}$$

Hence

$$m^{-1}(\mu) = \frac{k}{\mu + k}, \quad \begin{aligned} \mu(\mu+k) &= k \\ 0 &= \frac{k}{\mu+k} \end{aligned}$$

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 2.4.2 Method of moments for the Pareto distribution

Suppose X has density

$$p_\theta(x) = \theta(1+x)^{-(1+\theta)}, \quad x > 0, \quad \int_0^\infty x \cdot \theta (1+x)^{-(1+\theta)} dx$$

with respect to Lebesgue measure, and with $\theta \in \Theta \subset (0, \infty)$ (see Example 1.4.2).

Then, for $\theta > 1$

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

$$\int_0^\infty x \cdot \frac{\theta}{(1+x)^{1+\theta}} dx = x \left[\frac{\theta}{(1+x)^{1+\theta}} \right]_0^\infty = \left[\frac{x \cdot \theta \cdot (1+\theta)}{(1+x)^{1+\theta}} \right]_0^\infty = \left[\frac{\theta \cdot (1+\theta)}{(1+x)^{1+\theta}} \right]_0^\infty = \theta \cdot (1+\theta) \int_0^\infty \frac{1}{(1+x)^{1+\theta}} dx = \theta \cdot (1+\theta) \cdot \frac{1}{1+\theta} = \theta$$

$$= \left[\frac{\theta}{(1+\theta)^{1+\theta}} + -(1+\theta) \cdot \frac{\theta}{(1+\theta)^{2+\theta}} \right]_0^\infty$$

$$-\theta + (1+\theta) \cdot \theta = \cancel{\theta} + \theta^2 - \cancel{\theta}$$

20

CHAPTER 2. ESTIMATION

with inverse

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

The method of moments estimator would thus be

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

So not good for every case
 ! The function for the moment must exist and be possible injective and surjective over the entire parameter space.

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.

Notice: a measure ν is dominated by μ if $\nu \ll \mu$, that is if $\mu(A) = 0 \Rightarrow \nu(A) = 0$. A family $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$

2.4.3 Likelihood methods

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

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

Definition 2.4.1 The likelihood function (with data $\mathbf{X} = (X_1, \dots, X_n)$) is the function $L_{\mathbf{X}} : \Theta \rightarrow \mathbb{R}$ given by

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

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 footnotes in Sections 1.2 and 2.3). As we will see, different symbols are needed for the development of the theory.

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

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

Notice that you should take the log when the support makes sure you never end up with the untractable $\log(0)$.

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 $1/n$:

$$\hat{\theta} = \arg \max_{\vartheta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p_\vartheta(X_i). \quad (2.2)$$

despite the normalization the stays the same.

Replacing the average $\sum_{i=1}^n \log p_\vartheta(X_i)/n$ in (2.2) by its theoretical counterpart $E_\vartheta \log p_\vartheta(X)$ gives

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

which is indeed equal to the parameter θ we are trying to estimate:

Lemma 2.4.1 *We have*

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

Proof. By the inequality $\log x \leq x - 1$, $x > 0$, for all $\vartheta \in \Theta$

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

this is a constant; you can add without shifting the maximum

$$\begin{aligned} Q(\hat{\theta}_n) &= \arg \max_{\vartheta \in \Theta} E_{\hat{F}_n} \log p_\vartheta(x) \\ &= \cdots \frac{1}{n} \sum_{i=1}^n \log p_\vartheta(x) \\ &= \frac{1}{n} \log L_X(\vartheta) = \hat{\theta} \end{aligned}$$

*Since $F = F_\theta$, θ is true parameter
so that $Q(F) = Q(\theta)$*

$$\begin{aligned} Q(F) &= \arg \max_{\vartheta \in \Theta} E_\vartheta \log p_\vartheta(x) \\ &= \cdots E_\vartheta \log \left(\frac{p_\vartheta(x)}{p_\theta(x)} \right) \\ &\leq \cdots E_\vartheta \log \left(\frac{p_\theta(x)}{p_\theta(x)} \right) \\ &= \cdots \log E_\vartheta \left(\frac{p_\theta(x)}{p_\theta(x)} \right) \\ &= \cdots \log \int p_\theta(x) \frac{p_\theta(x)}{p_\theta(x)} dx \\ &= \cdots \log \int p_\theta(x) dx = \log(1) = 0 \end{aligned}$$

Example 2.4.3 MLE for the Pareto distribution

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

$$\log p_\vartheta(x) = \log \vartheta - (1 + \vartheta) \log(1 + x),$$

$$\frac{d}{d\vartheta} \log p_\vartheta(x) = \frac{1}{\vartheta} - \log(1 + x).$$

We put the derivative of the log-likelihood based on n observations 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 2.4.4 MLE for some location/scale models

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

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

parameters made explicitly

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

$$p_\theta(\cdot) = \frac{1}{\sigma} f_0\left(\frac{\cdot - \mu}{\sigma}\right). \quad \boxed{\frac{dF_0}{dx}}$$

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

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

first part of gradient

$$\frac{\partial \ln p_\theta(x)}{\partial \mu} = \frac{1}{\sigma^2} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)\right] = 0$$

$$\sum_{i=1}^n \frac{1}{\sigma^2} (x_i - \mu) = 0$$

$$n\mu = \sum_{i=1}^n x_i$$

$$\mu = \frac{\sum x_i}{n}$$

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

so

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

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

The MLE of μ resp. σ is now

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

$$\frac{\partial p_\theta(x)}{\partial \mu} = \frac{1}{2\sigma} \exp\left[-\frac{\sqrt{2}\sqrt{2}|x - \mu|}{\sigma}\right] = 0$$

$$= \text{sign}(x - \mu) = 0$$

This is minimized for the
 $\mu = \text{Median}.$

Example 2.4.5 An example where the MLE does not exist

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

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

Taking $\tilde{\mu} = X_1$ yields

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

because of this term.

Now, since for all $z \neq 0$

$$\lim_{\tilde{\sigma} \downarrow 0} \frac{1}{\tilde{\sigma}} \phi(z/\tilde{\sigma}) \xrightarrow{\text{goes to 0}} 0,$$

in goes infinite

we have

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

It follows that

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

Method of moments
works instead:

$$\hat{\mu} = \frac{1}{2}\hat{\theta}_1 + \frac{1}{2}\hat{\theta}_2^{-1},$$

$$\hat{M}_2 = \frac{1}{2} \left(1 + \hat{\theta}_2^{-1} \right) + \frac{1}{2} \left(\hat{\theta}_2^{-1} + \hat{\theta}_1^{-1} \right)$$

always positive

$$= \frac{1}{2} + \frac{1}{2}\hat{\theta}_2^{-1} + \hat{\theta}_1^{-1}$$

$$\hat{\theta}_1^{-1}, \hat{\theta}_2^{-1} \geq \hat{\theta}_2^{-1} - \hat{\theta}_1^{-1} = 2\hat{\theta}_2^{-1}$$

with $\hat{\theta}_2^{-1} \geq \frac{1}{2}(X_1 - \bar{X})^2$
no guarantee that $\hat{\theta}_2^{-1} > 0$

2.5 Asymptotic tests and confidence intervals based on the likelihood

This section is a look ahead for what is to come in Chapter 14. 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 14 that, under some regularity conditions,

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

Here, " $\xrightarrow{\mathcal{D}_{\theta}}$ " means convergence in distribution under \mathbb{P}_{θ} , and χ_p^2 denotes the Chi-squared distribution with p degrees of freedom.

We say that $Z(\mathbf{X}, \theta)$ is an asymptotic pivot: its asymptotic distribution does not depend on the unknown parameter θ . For the null-hypothesis

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

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

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

Example 2.5.1 Likelihood ratio for the normal distribution

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

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

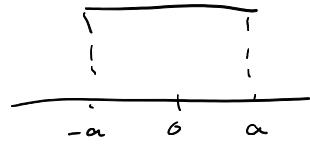
The random variable $\sqrt{n}(\bar{X} - \mu)$ is $\mathcal{N}(0, 1)$ -distributed under \mathbb{P}_{μ} . So its square, $n(\bar{X} - \mu)^2$, has a χ_1^2 -distribution. Thus, in this case the above test (confidence interval) is exact.

And not simply correct asymptotic in the exact case

As an exercise do
MLE location $\xrightarrow{\text{uniform}}$ Laplace dist. with $\text{var} = 1$
 $\xrightarrow{\text{normal}}$

$$\text{var} = 1 \quad \text{for Laplace} \Rightarrow f_0(x) = \frac{1}{\sqrt{2}} e^{-\sqrt{2}|x|}$$

$$f_0, \text{ vs form} \quad f_0(x) = \frac{1}{2a} \mathbb{1}_{[-a, a]}(x)$$



$$\text{with } a = \sqrt{3}$$

Chapter 3

Intermezzo: distribution theory

3.1 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

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

It follows that

$$P(B|A) = P(A|B) \frac{P(B)}{P(A)}, \quad \text{with } P(A \cap B) = P(A|B) \cdot P(B)$$

and that, for a partition $\{B_j\}$ ¹

$$P(A) = \sum_j P(A|B_j)P(B_j). \quad \left. \begin{array}{l} \text{the sum of} \\ \text{the probabilities} \\ \text{in each partition,} \\ \text{→ similar reasoning} \\ \text{to the marginal dist. reasoning} \end{array} \right\}$$

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

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

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

Thus, we have

$$f_Y(y|x) = f_X(x|y) \frac{f_Y(y)}{f_X(x)}, \quad (x, y) \in \mathbb{R}^{n+m}, \quad \left. \begin{array}{l} \text{same} \\ \text{reasoning} \\ \text{as above.} \end{array} \right\}$$

and

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

Definition 3.1.2 Let $g : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ be some function. The conditional expectation of $g(X, Y)$ given $Y = y$ is

$$E[g(X, Y)|Y = y] := \int \underbrace{f_X(x|y) g(x, y)}_{\text{probability density}} dx. \quad \left. \begin{array}{l} \text{integrating over} \\ \text{the other variable that is} \\ \text{not given by} \\ \text{conditional probability at } f(x|y). \end{array} \right\}$$

Note thus that

$$E(X) = \int x \cdot p(x) dx$$

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

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 3.1.1 (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

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

□

found on
the internet
early combination

3.2 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

Read:
 $\binom{n}{x,y,z}$, possible ways
 combinations to
 observe n cases. n is
 a feature like quality
 in a sample of size
 N.

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

↳ possible permutations/combinations of
 observing x, y, z times in a sample
 at n.

Here

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

It is called a *multinomial coefficient*.

Lemma 3.2.1 *The marginal distribution of X is the $\text{Binomial}(n, p_1)$ -distribution.*

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

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

some convergence reason as below, [Binomial series]

Definition 3.2.1 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

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

Here

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

Example 3.2.1 Histograms

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 $j = 1, \dots, k$,

$$\begin{aligned} \overline{p_j} &:= P(X \in (a_{j-1}, a_j]) = F(a_j) - F(a_{j-1}), \\ \overline{N_j} &:= \frac{\#\{X_i \in (a_{j-1}, a_j]\}}{n} = \hat{F}_n(a_j) - \hat{F}_n(a_{j-1}). \end{aligned}$$

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

3.3 The Poisson distribution

Definition 3.3.1 A random variable $X \in \{0, 1, \dots\}$ has the Poisson distribution with parameter $\lambda > 0$, if for all $x \in \{0, 1, \dots\}$

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

prob of x occurrences
in a time frame of length λ (avg of occurrences).

(see also Example 1.4.1).

average time for 1 occurrence
for a poisson dist. variable ($\Rightarrow \exp(\lambda) = \lambda \cdot e^{-\lambda}$)
with rate λ . If $x = 1$ you get the \square

immediate
to see
that it is
the taylor
expansion
of $(1+x)^a$
 $= \sum_{k=0}^{\infty} \binom{a}{k} x^k$

Lemma 3.3.1 Suppose X and Y are independent, and that X has the $\text{Poisson}(\lambda)$ -distribution, and Y the $\text{Poisson}(\mu)$ -distribution. Then $Z := X + Y$ has the $\text{Poisson}(\lambda + \mu)$ -distribution. *3-additive property.*

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

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

$$\begin{aligned} &\lambda^2 \cdot 0 \\ &+ \lambda \cdot \lambda^{z-1} \cdot \mu \\ &+ \epsilon \cdot \lambda' \mu^{z-1} \\ &z(\lambda^{z-1} \mu + \lambda^1 \mu^{z-1}) \end{aligned}$$

□

Lemma 3.3.2 Let X_1, \dots, X_n be independent, and (for $i = 1, \dots, n$), let X_i have the $\text{Poisson}(\lambda_i)$ -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 $\text{Poisson}(\lambda_+)$ -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

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

□

3.4 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

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

3.4. THE DISTRIBUTION OF THE MAXIMUM OF TWO RANDOM VARIABLES 29

Lemma 3.4.1 The distribution function of Z is F^2 . Moreover, Z has density

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

taking partial derivative

$$\begin{aligned} \text{cause } F(z) &= P(X_1 \leq z, X_2 \leq z) \quad \text{and if } X_1 \sim X_2 \text{ indep.} \\ F(z) &= F_1(z) \cdot F_2(z) \\ \text{as } F_1(\cdot) &= F_2(\cdot) \\ \text{by def} & \\ F(Z) &= F(\cdot)^2 \end{aligned}$$

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,

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

$$\begin{aligned} &\Pr[\min\{X_1, X_2\} \geq z] \\ &= 1 - \Pr[X_1 \leq z, X_2 \leq z] \\ &= 1 - F^2(z) \end{aligned}$$

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

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

□

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

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

you know
max 1 of the
two.

Note thus that this distribution has a jump of size $\frac{1}{2}$ at z .

look at this post
if really interested! I believe the answer is
hidden here.

<https://stats.stackexchange.com/questions/232085/conditional-distribution-of-uniform-random-variable-given-order-statistic>

$$\frac{f_{X_1|Z}}{f(z)} = \begin{cases} \frac{f(z)}{f(z)} = 1 & \text{if } x_1 \geq z \\ \frac{f(z)}{2f(z)} = \frac{1}{2} & \text{if } x_1 < z \end{cases}$$

$$\begin{aligned} f_{X_1|Z} &= \frac{f(x_1) \cdot \frac{1}{2}}{2f(z) \cdot f(z)} \\ &= \frac{\text{constant}}{\text{when integrating just}} \end{aligned}$$

Chapter 4

Sufficiency and exponential families

In this chapter, we denote the data by $X \in \mathcal{X}$. (In examples X is often replaced by $\mathbf{X} = (X_1, \dots, X_n)$ with X_1, \dots, X_n i.i.d. copies of X .) We assume X has distribution $P \in \{P_\theta : \theta \in \Theta\}$.

4.1 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).

Important Def.

Definition 4.1.1 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 .

If means that knowing the result of the given statistic the probability distribution of a Random Variable X is not affected by the current form of the parameteric solution θ . In this sense, the statistic is sufficient.

Example 4.1.1 Sufficiency and Bernoulli trials

Let $\mathbf{X} = (X_1, \dots, X_n)$ with X_1, \dots, X_n i.i.d. with 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 ,

$$\underbrace{\Pr_{\theta}^{(1-\theta)^{n-s}\theta^s} \frac{1}{\binom{n}{s}}}_{\Pr_{\theta}^{(1-\theta)^{n-s}\theta^s}} \sim \Pr_{\theta}\left(X_1 = x_1, \dots, X_n = x_n \mid S = s\right) = \frac{1}{\binom{n}{s}}, \quad \sum_{i=1}^n x_i = s.$$

Let S be sufficient (Given $s = s$). Let \mathbf{X}^* be a sample from the distribution of \mathbf{X} given $S = s$. Then \mathbf{X}^* has the same distribution as \mathbf{X} . $\Pr_{\theta}(x_i^* \in A_i) = \Pr_{\theta}(x_i \in A_i)$ $\Pr_{\theta}(x_i \in A_i) = \Pr_{\theta}(x_i \in A_i)$ $\Pr_{\theta}(x_i \in A_i) = \Pr_{\theta}(x_i \in A_i)$

Example 4.1.2 Sufficiency and the Poisson distribution

Let $\mathbf{X} := (X_1, \dots, X_n)$, with X_1, \dots, X_n i.i.d. and $\text{Poisson}(\theta)$ -distributed. Take

$$\begin{aligned} \Pr(A \cap B) &= \frac{\Pr(A \cap B)}{\Pr(B)} \\ &\downarrow A \subset B \\ &\Rightarrow = \frac{\Pr(A)}{\Pr(B)} \end{aligned}$$

Lemma 4.3.1
 $S = \sum_{i=1}^n X_i$. 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 $(p_1, \dots, p_n) = (\frac{1}{n}, \dots, \frac{1}{n})$: multinomial coefficient

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

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

Example 4.1.3 Sufficiency and the exponential distribution

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

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

Hence, S is sufficient for θ .

Example 4.1.4 Sufficiency of the order statistics

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 $(x_{q_1}, \dots, x_{q_n}) = s$,

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

Example 4.1.5 Sufficiency and the uniform distribution

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

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

not
the
derivation
of PDF

Any
permutation
is equally
likely.

Remember
 $f_{X+Y}(z) = f_X(a) \cdot f_Y(b)$

$$\frac{e^{itb} - e^{ita}}{it(b-a)} \rightarrow f_{Z \sim U[0,1]}(a, b)$$

With $X \wedge Y \sim U[a, b]$

$$\frac{e^{itb} - e^{ita}}{it(b-a)} \cdot \frac{e^{ita} - e^{ida}}{it(b-a)}$$

max
realization
 $Z = \max\{X, Y\}$

$$\frac{1}{\theta} \cdot \frac{1}{\theta} = \frac{1}{z}$$

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 .

4.2 Factorization Theorem of Neyman

Theorem 4.2.1 (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

$$p_\theta(x) = g_\theta(S(x)) \overset{\text{scoring constant}}{\sim} h(x), \quad \forall x, \theta$$

for some functions $g_\theta(\cdot) \geq 0$ and $h(\cdot) \geq 0$.

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

$$P_\theta(X = x | S = s) = \frac{P_\theta(X = x)}{Q_\theta(s)}, \quad S(x) = s. \quad (1)$$

\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 $S(x) = s$,

$$P_\theta(X = x) = P_\theta(X = x | S = s) Q_\theta(S = s) = h(x) g_\theta(s), \quad \square \text{ of the } \Rightarrow \text{ part.}$$

If and only if,

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

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

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

in θ above

does not depend on θ given sufficiency

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

*comes from inserting
in (4) above.*

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

which does not depend on θ .

□

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



Example 4.2.1 Sufficiency for the uniform distribution with unknown endpoint

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

$$p_\theta(x_1, \dots, x_n) = \frac{1}{\theta^n} \mathbb{1}\{0 \leq \min\{x_1, \dots, x_n\} \leq \max\{x_1, \dots, x_n\} \leq \theta\}$$

$S = X_{(n)}$ is sufficient.

with

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

and

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

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

Corollary 4.2.1 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 .

*to note $h(x)$ (for θ)
away cause it
does not depend on
 θ .*

4.3 Exponential families

Definition 4.3.1 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

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

sharpening
so that exp
pdf does not
 < 0 .
linear differentiable pdf.
makes sure that the pdf is integrated

*C: $\Theta \rightarrow \mathbb{R}$
T: $X \rightarrow \mathbb{R}^k$
d: $\Theta \rightarrow \mathbb{R}$
h: $X \rightarrow (0, \infty)$*
Note: $T(w) = (T_1(w), \dots, T_k(w))$ is sufficient.

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 $\mathbf{x} := (x_1, \dots, x_n)$),

$$p_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i) = \exp \left[\sum_{j=1}^k n c_j(\theta) T_j(\mathbf{x}) - n d(\theta) \right] \prod_{i=1}^n h(x_i),$$

Important

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 4.3.1 Poisson distribution

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 $d(\theta) = \theta$.

from this and note
(1) we see
immediately that if
 x_1, \dots, x_n i.i.d. Poisson(θ)
then $\sum x_i$ sufficient.

Example 4.3.2 Binomial distribution

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 4.3.3 Negative binomial distribution

If X has the Negative Binomial(m, θ)-distribution with m known we have

$$\begin{aligned} p_\theta(x) &= \frac{\Gamma(x+m)}{\Gamma(m)x!} \theta^m (1-\theta)^x \Rightarrow \text{inverted } x \wedge m. \text{ How is it defined for number of successes?} \\ &= \frac{\Gamma(x+m)}{\Gamma(m)x!} \exp[x \log(1-\theta) + m \log(\theta)]. \end{aligned}$$

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

Example 4.3.4 Gamma distribution with 1 parameter

Let X have the Gamma(m, θ)-distribution with m known. Then

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

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

Example 4.3.5 Gamma distribution with two parameters

Let X have the $\text{Gamma}(m, \lambda)$ -distribution, and let $\theta = (m, \lambda)$. Then

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

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

Example 4.3.6 Normal distribution

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 $d(\theta) = \mu^2/(2\sigma^2) + \log(\sigma)$.

4.4 Intermezzo: the mean and covariance matrix of a random vector

Let $Z \in \mathbb{R}^k$ be a random vector. Then its mean (if it exists) is defined as the vector consisting of the means of each entry of Z :

$$EZ = \begin{pmatrix} EZ_1 \\ \vdots \\ EZ_k \end{pmatrix}.$$

The covariance matrix of X (if it exists) is defined as the symmetric $k \times k$ matrix Σ containing the (co)variances between each pair of entries in X :

$$\Sigma := EZZ' - EZEZ' = \begin{pmatrix} \sigma_1^2 & \sigma_{1,2} & \cdots & \sigma_{1,k} \\ \sigma_{1,2} & \sigma_2^2 & \cdots & \sigma_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1,k} & \sigma_{2,k} & \cdots & \sigma_k^2 \end{pmatrix}.$$

Here, Z' denotes the transpose¹ of the vector Z , $\sigma_j^2 := \text{var}(Z_j)$ ($j = 1, \dots, k$) and for all $j_1 \neq j_2$ $\sigma_{j_1, j_2} = \text{cov}(Z_{j_1}, Z_{j_2})$. We will often write the covariance matrix Σ as $\text{Cov}(Z)$.

¹We alternatively write the transpose of a vector, say v , as v^T .

4.5 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

$$\dot{d}(\theta) := \frac{\partial}{\partial \theta} d(\theta)$$

denote the vector of first derivatives. Let

$$\ddot{d}(\theta) := \frac{\partial^2}{\partial \theta \partial \theta'} d(\theta) = \left(\frac{\partial^2}{\partial \theta_{j_1} \partial \theta_{j_2}} d(\theta) \right)$$

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

$$\text{Cov}_\theta(T(X)) := E_\theta T(X) T'(X) - E_\theta T(X) E_\theta T'(X).$$

$E_\theta = E(x^2) - E(x)^2$

Lemma 4.5.1 We have (under regularity)

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

With
canonical
form.

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

$$\begin{aligned}
 d(\theta) &= \frac{\partial}{\partial \theta} \log \left(\int \exp[\theta' T(x)] h(x) d\nu(x) \right) \\
 &\stackrel{\text{assume possible}}{=} \frac{\int \exp[\theta' T(x)] T(x) h(x) d\nu(x)}{\int \exp[\theta' T(x)] h(x) d\nu(x)} \\
 &= \int \exp[\theta' T(x) - d(\theta)] T(x) h(x) d\nu(x) \\
 &= \int p_\theta(x) T(x) d\nu(x) = E_\theta T(X),
 \end{aligned}$$

$\frac{d \log(\cdot)}{dx} = \frac{f'(x)}{f(x)}$
 $\int e^{[\theta' T(x)]} \cdot T(x) h(x) d\nu(x) \cdot (e^{[\theta' T(x)]})^{-1}$
 $= \int e^{(\theta' T(x))} \cdot T(x) h(x) d\nu(x) \cdot -\log(e^{\theta' T(x)} \cdot h(x) d\nu(x))$
 $= d(\theta), \text{ by definition}$

and (omitting the integration x variable to shorten the expressions)

$$\begin{aligned}
 \ddot{d}(\theta) &= \frac{\int \exp[\theta' T] TT' h d\nu}{\int \exp[\theta' T] h d\nu} \\
 &- \frac{\left(\int \exp[\theta' T] Th d\nu \right) \left(\int \exp[\theta' T] Th d\nu \right)^T}{\left(\int \exp[\theta' T] h d\nu \right)^2} \\
 &= \int \exp[\theta' T - d(\theta)] TT' h d\nu \\
 &- \left(\int \exp[\theta' T - d(\theta)] Th d\nu \right) \times \left(\int \exp[\theta' T - d(\theta)] T' h d\nu \right) \\
 &= \int TT' p_\theta d\nu - \left(\int p_\theta T d\nu \right) \left(\int p_\theta T' d\nu \right) \\
 &= E_\theta T(X) T'(X) - \left(E_\theta T(X) \right) \left(E_\theta T'(X) \right) \\
 &= \text{Cov}_\theta(T(X)).
 \end{aligned}$$

quotienten regel:
 $f(x) = \frac{v(x)}{u(x)}$
 $f'(x) = \frac{u'(x) \cdot v(x) - v'(x) \cdot u(x)}{u(x)^2}$

□

4.6 Reparametrizing in the one-dimensional case

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

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

\downarrow
 $d(c(\theta))$ is the form of the canonical family above.

to get

$$\tilde{p}_\gamma(x) = \exp[\gamma T(x) - \tilde{d}_0(\gamma)] h(x),$$

where when c is one-to-one invertible

$$\begin{aligned} d_0(\gamma) &= d(c^{-1}(\gamma)). \\ \frac{d_0(\gamma)}{d\theta} &= \dot{d}(\gamma) \cdot \frac{d(\gamma)}{d\theta} = \dot{d}(\gamma) \cdot \frac{d(\theta)}{d\theta} = \dot{d}(c^{-1}(\gamma)) \end{aligned}$$

Umkehrregel
der Ableitung.

It follows that

$$E_\theta T(X) = \dot{d}_0(\gamma) = \frac{\dot{d}(c^{-1}(\gamma))}{\dot{c}(c^{-1}(\gamma))} = \frac{\dot{d}(\theta)}{\dot{c}(\theta)}, \quad (4.1)$$

by the derivation at before

and

$$\begin{aligned} \text{var}_\theta(T(X)) &= \ddot{d}_0(\gamma) = \frac{\ddot{d}(c^{-1}(\gamma))}{[\dot{c}(c^{-1}(\gamma))]^2} = \frac{\dot{d}(c^{-1}(\gamma)) \ddot{c}(c^{-1}(\gamma))}{[\dot{c}(c^{-1}(\gamma))]^3} \\ &= \frac{\ddot{d}(\theta)}{[\dot{c}(\theta)]^2} - \frac{\dot{d}(\theta) \ddot{c}(\theta)}{[\dot{c}(\theta)]^3} \\ &= \frac{1}{[\dot{c}(\theta)]^2} \left(\ddot{d}(\theta) - \frac{\dot{d}(\theta)}{\dot{c}(\theta)} \ddot{c}(\theta) \right). \end{aligned} \quad (4.2)$$

from innerhalb der
Bilanzierungsregel:
 $p(x) = \frac{U(x)}{V(x)}$
 $L'(x) = \frac{U'(x) \cdot V(x) - U(x) \cdot V'(x)}{V(x)^2}$

this cancels out with the $U(x)$ term here

4.7 Score function and Fisher information

Consider an arbitrary (but regular) family of densities $\{p_\theta : \theta \in \Theta\}$, with (again for simplicity) $\Theta \subset \mathbb{R}$.

Definition 4.7.1 The score function is

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

The Fisher information for estimating θ is

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

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

See also Chapters 5 and 13.

Lemma 4.7.1 We have (under regularity)

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

and

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

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:

$$\begin{aligned} E_\theta s_\theta(X) &= \int s_\theta(x)p_\theta(x)d\nu(x) \\ &= \int \frac{d \log p_\theta(x)}{d\theta} p_\theta(x)d\nu(x) = \int \frac{\dot{p}_\theta(x)}{p_\theta(x)} p_\theta(x)d\nu(x) \\ &= \int \dot{p}_\theta(x)d\nu(x) = \left(\frac{d}{d\theta} \right) \int p_\theta(x)d\nu(x) = \frac{d}{d\theta} \mathbf{1} = 0, \end{aligned}$$

and

$$\begin{aligned} E_\theta \dot{s}_\theta(X) &= E_\theta \left[\frac{\ddot{p}_\theta(X)}{p_\theta(X)} - \left(\frac{\dot{p}_\theta(X)}{p_\theta(X)} \right)^2 \right] \\ &= E_\theta \left[\frac{\ddot{p}_\theta(X)}{p_\theta(X)} \right] + E_\theta \dot{s}_\theta^2(X) = \text{var}_\theta s_\theta(X). \end{aligned}$$

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

$$\begin{aligned} \text{var}(x_i) &= E[(x_i - \bar{x})^2] \\ &= E[x_i^2 - 2x_i\bar{x} + \bar{x}^2] \\ &= E(x_i^2) - 2\bar{x}^2 + \bar{x}^2 \\ &= E(x_i^2) - E(x_i)^2 \end{aligned}$$

$$\begin{aligned} E_\theta \left[\frac{\ddot{p}_\theta(X)}{p_\theta(X)} \right] &= \int \frac{d^2}{d\theta^2} p_\theta(x)d\nu(x) \\ &= \frac{d^2}{d\theta^2} \int p_\theta(x)d\nu(x) \\ &= \frac{d^2}{d\theta^2} \mathbf{1} = 0. \end{aligned}$$

□

4.8 Score function for exponential families

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

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

$$\frac{\log(\cdot)}{d\theta} \Leftrightarrow \frac{\text{dot } d}{\text{Score function}}$$

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

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

which re-establishes (4.1). 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

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

$$\begin{aligned} \text{var}_\theta(s_\theta(x)) &= \text{var}_\theta(\dot{c}(\theta)T(x) - \dot{d}(\theta)) \\ &= \text{var}_\theta(\dot{c}(\theta)T(x)) + \text{var}_\theta(-\dot{d}(\theta)) \\ &= \text{var}_\theta(\dot{c}(\theta)T(x)) \\ &\quad \circ \dot{c}(\theta) \text{ var}(T(x)) \end{aligned}$$

4.8. SCORE FUNCTION FOR EXPONENTIAL FAMILIES

41

which re-establishes (4.2). In addition, it follows that

$$I(\theta) = \ddot{d}(\theta) - \frac{\dot{d}(\theta)}{\dot{c}(\theta)} \ddot{c}(\theta). \quad \left\{ \begin{array}{l} \text{for one-dimensional} \\ \text{exponential families.} \end{array} \right.$$

big mass
for nothing

$$\ddot{d}_0(\gamma) = \text{var}_{\theta}(T(X)) \quad \text{and} \quad \text{var}_{\theta}(T(X)) = \frac{1}{\dot{c}(\theta)^2}$$

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

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

this part follows by definition

Example 4.8.1 Bernoulli-distribution in canonical form

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

$E_{\theta}(X) = \dot{d}(\theta)$
where
reparametrization

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

We reparametrize:

$$\gamma := c(\theta) = \log\left(\frac{\theta}{1-\theta}\right), \quad \Rightarrow \text{you reparametrize the } \theta \text{ in the canonical form above}$$

which is called the log-odds ratio. Inverting gives

$$\theta = \frac{e^{\gamma}}{1+e^{\gamma}}, \quad \begin{aligned} & -\log\left(1 - \frac{e^{-\gamma}}{1+e^{-\gamma}}\right) \\ &= -\log\left(\frac{1+e^{-\gamma}}{1+e^{-\gamma}}\right) \\ &= -\log\left(\frac{1}{1+e^{-\gamma}}\right) + \log(1+e^{-\gamma}) \\ d(\theta) = -\log(1-\theta) &= \log\left(1+e^{\gamma}\right) := d_0(\gamma). \end{aligned}$$

$$\log\left(\int \exp(x \cdot \log\left(\frac{\theta}{1-\theta}\right)) dv\right)$$

$$\log\left(\sum_{x=0}^1 \exp(x \cdot \log\left(\frac{\theta}{1-\theta}\right))\right) \stackrel{\text{def}}{=} \log\left(\log\frac{\theta}{1-\theta}\right)$$

and hence

by definition exponential family

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

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

$$\begin{aligned} s_{\theta}(x) &= \frac{d}{d\theta} \left[x \log\left(\frac{\theta}{1-\theta}\right) + \log(1-\theta) \right] \\ &= \frac{x}{\theta(1-\theta)} - \frac{1}{1-\theta}. \quad \left\{ \begin{array}{l} = \frac{x - \theta}{\theta(1-\theta)} \\ E_{\theta}s_{\theta}^2 = \frac{1}{\theta(1-\theta)} \left(\frac{(x-\theta)^2}{\theta(1-\theta)} \right) \end{array} \right. \end{aligned}$$

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

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

4.9 Minimal sufficiency

Free proportional given two data x and \tilde{x}

Definition 4.9.1 We say that two likelihoods $L_x(\theta)$ and $L_{\tilde{x}}(\theta)$ are proportional at (x, \tilde{x}) , if

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

for some constant $c(x, \tilde{x})$.

We write this as

$$L_x(\theta) \propto L_{\tilde{x}}(\theta).$$

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



sufficient statistic
stays the same
 \rightarrow you can see this by the factorization theorem

Example 4.9.1 Minimal sufficiency for the normal distribution

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

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

So

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

which equals,

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

notice these like this logs. is the exact result that matches the definition above.

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

Example 4.9.2 Minimal sufficiency for the Laplace distribution

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

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

so

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

which equals

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

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

Chapter 5

Bias, variance and the Cramér Rao lower bound

5.1 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 is $\gamma := g(\theta)$, with $g : \Theta \rightarrow \mathbb{R}$. Except for the last section in this chapter, the parameter γ is assumed to be one-dimensional.

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

Definition 5.1.1 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 .

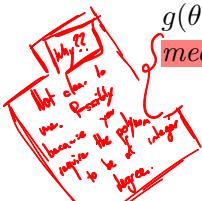
Example 5.1.1 Unbiased estimators in the Binomial case

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 $\theta/(1-\theta)$.

case otherwise in expectation
higher terms will remain and
some persistent bias will occur.



Example 5.1.2 Unbiased estimators in the Poisson caseLet $X \sim \text{Poisson}(\theta)$. Then

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

same as taking the integral in the discrete case

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

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

As another example, suppose the parameter of interest is

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

An unbiased estimator is

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

Understand

\hookrightarrow Prob. due to convergence of the exponent.

This estimator does not make sense at all!

Example 5.1.3 Unbiased estimator of the varianceLet 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

Idea:

$$\begin{aligned} &\text{use the trick} \\ &E((X_i - \mu) - (\bar{X} - \mu))^2 \\ &= E(X_i - \bar{X})^2 \end{aligned}$$

$$S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \text{by Jensen's inequality} \quad E\sqrt{S^2} \leq \sqrt{E(S^2)}$$

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.

!!

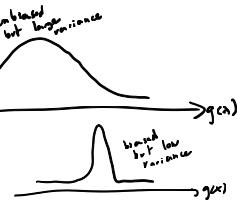
$$\begin{aligned} &E(\sigma^2 - 2(X - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2) \\ &= \sigma^2(n-1) \end{aligned}$$

5.2 UMVU estimators**Definition 5.2.1** The mean square error of T is

$$\text{MSE}_\theta(T) := E_\theta \left(T - g(\theta) \right)^2.$$

Lemma 5.2.1 We have the following decomposition for the mean square error:

$$\text{MSE}_\theta(T) = \text{bias}_\theta^2(T) + \text{var}_\theta(T).$$



$$E_\theta(T - g(\theta))^2 = E_\theta(T^2) - 2g(\theta)E_\theta(T) + g(\theta)^2$$

Proof. Write $E_\theta := q(\theta)$. Then

$$\begin{aligned} E_\theta(T - g(\theta))^2 &= \underbrace{E_\theta(T^2)}_{=\text{var}_\theta(T)} - 2\underbrace{g(\theta)E_\theta(T)}_{=\text{bias}_\theta^2(T)} + g(\theta)^2 \\ &\quad + 2(g(\theta) - q(\theta)) \underbrace{E_\theta(T - q(\theta))}_{=0}. \end{aligned}$$

$\rightarrow \text{As } E_\theta(T) = q(\theta) \quad \square$

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 5.2.1 Estimators in case of the normal distribution

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 $T := a\bar{X}$, where $0 \leq a \leq 1$. Then the bias is decreasing in a , but the variance is increasing in a :

$$\text{MSE}_\theta(T) = E_\theta(T - \mu)^2 = (1-a)^2\mu^2 + a^2\sigma^2/n.$$

↑ typical
bias-variance trade-off.

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

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

$T_{\text{opt}} = \left(\frac{\bar{X}^2}{\frac{\sigma^2}{n} + \bar{X}^2} \right)$.
see Stein's Estimator/
p.122.

However, $a_{\text{opt}}\bar{X}$ 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

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

Scripts omits
that holds just
based on the assumption
data are sampled from a Normal.

Then of course it holds that, the normalized square data follow asymptotically $\sim \chi^2$
 $\left(\frac{\bar{X}(X_i - \bar{X})}{\sigma}\right)^2 \sim \chi_{n-1}^2$

CHAPTER 5. BIAS, VARIANCE AND THE CRAMÉR RAO LOWER BOUND

$$\text{Var}\left(\frac{(n-1)S^2}{\sigma^2}\right) = 2(n-1)$$

$$\frac{(n-1)^2}{n} \text{Var}(S^2) = 2(n-1)$$

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

$$\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} \cdot \sigma^4$$

Thus¹

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

It follows that



$$\text{MSE}_\theta(S^2) = E_\theta \left(S^2 - \sigma^2 \right)^2 = \text{var}(S^2) = \frac{2(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{n-1}.$$

Moreover

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

so that

$$\begin{aligned} \text{MSE}_\theta(\hat{\sigma}^2) &= E_\theta \left(\hat{\sigma}^2 - \sigma^2 \right)^2 \\ &= \text{bias}_\theta(\hat{\sigma}^2)^2 + \text{var}_\theta(\hat{\sigma}^2) \\ &= \frac{\sigma^4}{n^2} + \frac{2(n-1)\sigma^4}{n^2} = \frac{(2n-1)\sigma^4}{n^2}. \end{aligned}$$

$$\begin{aligned} &\left(\frac{n-1}{n} \sigma^2 - \sigma^2 \right) \\ &= -\frac{1}{n} \sigma^2 \end{aligned}$$

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 5.2.2 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: by the iterated expectations lemma

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

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 :

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

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

Exercise compute minimum bias

“ $E(T - \sigma^2)$.”

solution

$\text{bias} = \boxed{1/n}$

So not only MLE estimator better in terms of MSE but also restriction. Being objective you should use the formula above.

$$\begin{aligned} \text{Intermezzo: } \text{Var}(Y) &= E[\text{Var}(Y|S)] + \text{Var}[E(Y|S)] \\ \text{Proof: } E[\text{Var}(Y|S)] &= E[(E(Y|S) - E(Y|S))^2] \\ &= E[(E(Y|S) - E(Y))^2] \\ \text{Var}[E(Y|S)] &= E[(E(Y|S))^2] - (E[E(Y|S)])^2 \\ &= E[(Y|S)^2] - (E[Y|S])^2 \end{aligned}$$

when adding

$= E(Y^2) - (E(Y))^2$

Lemma 5.2.2 Let Y and Z be two random variables. Then

$$\text{var}(Y) = \underbrace{\text{var}(E(Y|Z))}_{\substack{\text{variance of} \\ \text{conditional expectation}}} + \underbrace{\text{Evar}(Y|Z)}_{\substack{\text{expected conditional variance}}}. \quad \| !$$

Proof. It holds that

$$\begin{aligned} \text{var}(E(Y|Z)) &= E\left[E(Y|Z)\right]^2 - \underbrace{\left[E(E(Y|Z))\right]^2}_{\substack{\text{by def} \\ \text{iterated expectations}}}^2 \\ &= E[E(Y|Z)]^2 - [EY]^2, \end{aligned}$$

and

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

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

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

□

5.3 The Lehmann-Scheffé Lemma

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

$$E_\theta h(S) = 0 \forall \theta \Rightarrow h(S) = 0, \quad P_\theta - \text{a.s.}, \quad \text{almost surely w.r.t. the probability measure } P_\theta.$$

Here, h is a function of S not depending on θ .

There \hookrightarrow There \Rightarrow positive probability mass.

Lemma 5.3.1 (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 $\text{var}_\theta(T^*) \leq \text{var}_\theta(T) \forall \theta$. If $T'(S)$ is another unbiased estimator of $g(\theta)$, we have

$$E_\theta(T(S) - T'(S)) = 0, \forall \theta. \quad \left\{ \begin{array}{l} \text{expectation is a linear operator. Moreover both} \\ \text{unbiased; that is in expectation equal so that the difference equals 0.} \end{array} \right.$$

Because S is complete, this implies

$$\left\| \begin{array}{l} \text{and the difference above} \\ \text{might be interpreted} \\ \text{as the } h(S) \text{ of} \\ \text{the definition above.} \end{array} \right. \quad T^* = T', \quad P_\theta - \text{a.s.}$$

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 5.3.1 UMVU estimator in the Poisson case

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

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

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 ,

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

$\underbrace{P_\theta(k)}$

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

$$P(X=0) = \frac{e^{-\theta}}{1}$$

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

"
 $e^{-\theta}$

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. Let us now calculate T^* . First,

$$\begin{aligned} P(X=0, S=s) \\ \hline P(S=s) \end{aligned}$$

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

symmetric
↓ $\alpha - \ell$

Hence

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

is UMVU.

Example 5.3.2 UMVU estimation for uniform distribution

Let X_1, \dots, X_n be i.i.d. Uniform[0, θ]-distributed, and $g(\theta) := \theta$. We know

$$F(x_1) \xrightarrow{h} \left(\frac{s}{\theta}\right)^n$$

5.4. COMPLETENESS FOR EXPONENTIAL FAMILIES

49

that $S := \max\{X_1, \dots, X_n\}$ is sufficient (see Example 4.2.1). 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.

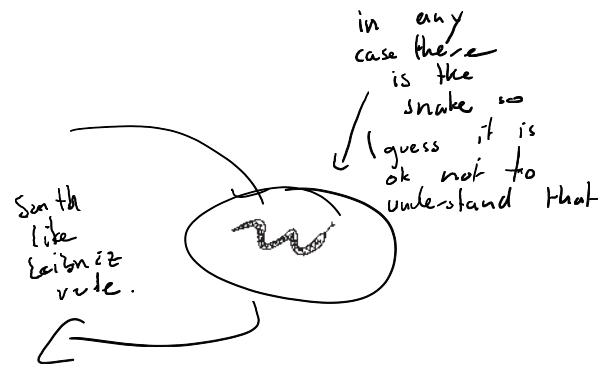
We have

$$E_\theta S = \int_0^\theta s \frac{ns^{n-1}}{\theta^n} ds = \frac{n}{n+1}\theta. \quad = \int_0^\theta \frac{ns^n}{\theta^n} ds = \left[\frac{n}{n+1} s \frac{\theta^{n+1}}{\theta^n} \right]_0^\theta = \frac{n}{n+1}\theta$$

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.



5.4 Completeness for exponential families

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 5.4.1 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 5.4.1 Completeness for Gamma distribution

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

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

It follows that

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

2 dim

is sufficient and complete.

Example 5.4.2 Completeness for two normal samples

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

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

must respect
this dimension

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 sufficient and complete statistic exist?

5.5 The Cramér 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

it holds that

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

derivative
 must move towards $s_\theta(x)$ in limit.
 must converge to a constant $(X) \in \mathbb{R}$.

Definition 5.5.1 If I and II hold, we call s_θ the score function, and $I(\theta)$ the Fisher information.

Comparing the above with Definition 4.7.1, we see that they coincide if $\theta \mapsto p_\theta$ is differentiable and “regularity conditions” hold. Recall also that in Lemma 4.7.1 we assumed unspecified “regularity conditions”. We now present a rigorous proof of the first part of this lemma.

Lemma 5.5.1 Assume Conditions I and II. Then

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

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,

$$\begin{aligned}
 |E_\theta s_\theta(X)|^2 &= \left| E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{hp_\theta(X)} - s_\theta(X) \right) \right|^2 \\
 &\leq \underbrace{\left[E_\theta \left(\frac{p_{\theta+h}(X) - p_\theta(X)}{hp_\theta(X)} - s_\theta(X) \right) \right]^2}_{\substack{\rightarrow 0 \\ \text{by condition II}}} \quad \text{Jensen's inequality}
 \end{aligned}$$

□

Note Thus $I(\theta) = \text{var}_\theta(s_\theta(X))$. *[Follows here and in line with what is shown in chapter 4]*

Remark As already noted, if $p_\theta(x)$ is differentiable for all x , we can take (under regularity conditions)

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

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

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

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 5.5.1 (The Cramér-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,

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

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

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

*[squared becomes Schwarz]
As proved above.
 $\int T^2 p_\theta \cdot \left(\frac{p_{\theta+h} - p_\theta}{h p_\theta} - s_\theta \right)^2 p_\theta = 0$*

*[let T be a possibly biased estimate of $g(\theta)$.
write $q(\theta) := E_\theta T(X)$
 $\text{MSE}_\theta(T) = \text{Var}_\theta(T) + \text{bias}^2(T)$
 $\geq \frac{(\dot{g}(\theta))^2}{I(\theta)} + \frac{(q(\theta) - g(\theta))^2}{(q(\theta) - g(\theta))^2}$
↳ then you can minimize this one to get to your desired result.]*

*[$\text{Cov}(T, s_\theta(X))$ as $E(s_\theta)$ has mean 0!
↳ As $\text{Cov}(T, s_\theta(X)) = E(T s_\theta) - E[T] E[s_\theta]$*

[$E(g(\theta) \cdot A) = E(g(\theta)) \cdot E(A)$ can take $g(\theta)$ as it does not depend on A . No exp. needed.]

as $h \rightarrow 0$. This is because by the Cauchy-Schwarz inequality

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

if s_θ var goes to 0 then so does non approx.

Thus,

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

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

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

□

Definition 5.5.2 We call $\dot{g}^2(\theta)/I(\theta)$, $\theta \in \Theta$, the Cramer Rao lower bound (CRLB) (for estimating $g(\theta)$).

Example 5.5.1 CRLB for exponential case

Let X_1, \dots, X_n be i.i.d. $\text{Exponential}(\theta)$, $\theta > 0$. 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,

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

$$\begin{aligned} \text{E}(x) &= 1/\theta & \text{E}(s_\theta(x)) &= 0 \\ \text{I}(\theta) &= \text{var}_\theta(X) = \frac{1}{\theta^2}. \end{aligned}$$

The CRLB for n observations is thus

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

In other words, T reaches the CRLB.

Example 5.5.2 CRLB for Poisson case 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

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

$$g(\theta) = e^{-\theta} \quad g'(\theta) = -e^{-\theta}$$

$$\text{CRLB} > \frac{\theta e^{-\theta}}{n}$$

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

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

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

5.6 CRLB and exponential families

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 5.6.1 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 Cramér 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 5.5, when T reaches the CRLB, we must have

$$\text{var}_\theta(T) = \frac{|\text{cov}_\theta(T, s_\theta(X))|^2}{\text{var}_\theta(s_\theta(X))}, \quad \Leftrightarrow \quad 1 = \frac{|\text{cov}(T, s_\theta)|^2}{\text{var}(s_\theta) \cdot \text{var}(T)} \quad \begin{matrix} \text{if } \text{cov}(T, s_\theta) \\ \text{is } \sqrt{\text{var}(T) \cdot \text{var}(s_\theta)} := \text{corr} \end{matrix}$$

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

perfect linear relation.

$$s_\theta(X) = a(\theta)T(X) - b(\theta). \quad (5.1)$$

if you think this exactly equals $\frac{p_\theta}{\dot{p}_\theta}$.

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 (5.1) 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 $g(\theta) = b(\theta)/a(\theta)$. \square

5.7 Higher-dimensional extensions

Expectations and covariance matrices of random vectors

Let $Z \in \mathbb{R}^k$ be a k -dimensional random vector. Then EZ is a k -dimensional vector, and

$$\Sigma := \text{Cov}(Z) := EZZ' - (EZ)(EZ') \quad \left\{ \begin{array}{l} \text{variance} \\ = E[Z^2] - E[Z]^2 \end{array} \right. \quad \begin{matrix} \text{as } \det \text{Cov.} \end{matrix}$$

For positive semi-definite Matrix it holds.

$$V > 0 \Rightarrow \text{pos-def}$$

$$V = Q \Lambda Q^T$$

$$Q^T Q = I$$

$Q = [q_1, \dots, q_k]$ square

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$ eigenval

$$V^{1/2} = Q \Lambda^{1/2} Q^T$$

Then

$$V = V^{1/2} V^{1/2}$$

CHAPTER 5. BIAS, VARIANCE AND THE CRAMÉR RAO LOWER BOUND

is a $k \times k$ matrix containing all variances (on the diagonal) and covariances (off-diagonal). Note that Σ is positive semi-definite: for any vector $a \in \mathbb{R}^k$, we have

$$\text{var}(a' Z) = a' \Sigma a \geq 0.$$

$$\begin{aligned} \text{proof: } & \text{Var}(a^T Z) = E[a^T (Z - \mu)]^2 \\ & = E[a^T (Z - \mu)(Z - \mu)^T] \\ & = a^T \text{Cov } a \end{aligned}$$

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

$$\max_{a \in \mathbb{R}^p} \frac{|a' c|^2}{a' V a} = c' V^{-1} c.$$

Proof. Write $V = W^2$, and $b := Wa$, $d := W^{-1}c$. Then $a' V a = b'b = \|b\|^2$ and $a' c = b'd$. By Cauchy-Schwarz

$$\max_{b \in \mathbb{R}^p} \frac{|b'd|^2}{\|b\|^2} = \|d\|^2 = d'd = c' V^{-1} c.$$

□

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 $\Theta \subset \mathbb{R}^k$. Let

$$g : \Theta \rightarrow \mathbb{R},$$

be a given function. Denote the vector of partial derivatives as

$$\dot{g}(\theta) := \begin{pmatrix} \partial g(\theta)/\partial \theta_1 \\ \vdots \\ \partial g(\theta)/\partial \theta_k \end{pmatrix}.$$

The score vector is defined as

$$s_\theta(\cdot) := \begin{pmatrix} \partial \log p_\theta/\partial \theta_1 \\ \vdots \\ \partial \log p_\theta/\partial \theta_k \end{pmatrix}.$$

The Fisher information matrix is

$$I(\theta) = E_\theta s_\theta(X) s'_\theta(X) = \text{Cov}_\theta(s_\theta(X)).$$

↳ As expectation = 0

possible to decompose positive definite matrix in the product of two positive definite matrices

Theorem 5.7.1 Let T be an unbiased estimator of $g(\theta)$. Then, under regularity conditions,

$$\text{var}_\theta(T) \geq \dot{g}(\theta)' I(\theta)^{-1} \dot{g}(\theta).$$



Proof. As in the one-dimensional case, one can show that, for $j = 1, \dots, k$,

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

Hence, for all $a \in \mathbb{R}^k$,

$$\begin{aligned} |a' \dot{g}(\theta)|^2 &= |\text{cov}_\theta(T, a' s_\theta(X))|^2 && \text{by definition / Cauchy-Schwarz.} \\ &\leq \text{var}_\theta(T) \text{var}_\theta(a' s_\theta(X)) && \text{by def. of var.} \\ &= \text{var}_\theta(T) a' I(\theta) a. && \text{F} \end{aligned}$$

Combining this with the auxiliary lemma gives

$$\text{var}_\theta(T) \geq \max_{a \in \mathbb{R}^k} \frac{|a' \dot{g}(\theta)|^2}{a' I(\theta) a} = \dot{g}'(\theta) I(\theta)^{-1} \dot{g}(\theta). \quad \text{lemma part.}$$

□

Corollary 5.7.1 As a consequence, one obtains a lower bound for unbiased estimators of higher-dimensional parameters of interest. As example, let $g(\theta) := \theta = (\theta_1, \dots, \theta_k)'$, and suppose that $T \in \mathbb{R}^k$ is an unbiased estimator of θ : $E_\theta T = \theta \forall \theta$. Then, for all $a \in \mathbb{R}^k$, $a'T$ is an unbiased estimator of $a'\theta$. Since $a'\theta$ has derivative a , the CRLB gives

$$\text{var}_\theta(a'T) \geq a'I(\theta)^{-1}a.$$

But

$$\text{var}_\theta(a'T) = a' \text{Cov}_\theta(T) a.$$

So for all a ,

$$a' \text{Cov}_\theta(T) a \geq a'I(\theta)^{-1}a,$$

in other words, $\text{Cov}_\theta(T) \geq I(\theta)^{-1}$, that is, $\text{Cov}_\theta(T) - I(\theta)^{-1}$ is positive semi-definite.

Chapter 6

Tests and confidence intervals

6.1 Intermezzo: quantile functions

Let F be a distribution function on \mathbb{R} . Then F is *cadlag* (continue à droite, limite à gauche). Define the quantile functions

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

the highest
 x for which
 F(x) ≤ u
 holds

and

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

It holds that

$$\boxed{F(q_{\inf}^F(u)) \geq u}$$

and, for all $h > 0$,

$$\boxed{F(q_{\sup}^F(u) - h) \leq u.}$$

Hence

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

6.2 How to construct tests

Consider a model class

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

Moreover, consider a space Γ , and a map

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

We think of γ as the parameter of interest.

Definition 6.2.1 Let $\gamma_0 \in \Gamma$ and $\alpha \in [0, 1]$ be given. A (non-randomized) test at level α for the hypothesis

$$H_0 : \gamma = \gamma_0$$

(is a statistic $\phi(X, \gamma_0) \in \{0, 1\}$ such that $P_\theta(\phi(X, \gamma_0) = 1) \leq \alpha$ for all $\theta \in \{\vartheta : g(\vartheta) = \gamma_0\}$)

does not depend on P & P is DPP in

↓
 this is why non-randomized.
 if you randomize, expectation
 definition is better.

that is
 given that
 the parameter
 is γ_0 the probability
 must be
 $P(X, \gamma_0)$
 always.

We often omit the dependence of ϕ on γ_0 , i.e., we write $\phi(X) := \phi(X, \gamma_0)$.

Note Typically a test ϕ is based on a test statistic T , i.e. it is of the form

$$\phi(X) = \begin{cases} 1 & \text{if } T(X) > c \\ 0 & \text{else} \end{cases} \quad \begin{array}{l} \hookrightarrow \text{so it is a statistic} \\ \text{based on a statistic.} \end{array}$$

The constant c is called the *critical value*.

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

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

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 $\gamma_0 = g(\theta_0)$:

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

And notice our pivot does not depend on θ such that the test statistics are valid across methods.

Asymptotic pivot 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 $\theta \in \Theta$,

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

at all continuity points of G , where the limit G does not depend on θ .

Example 6.2.1 The location model

As example, consider again the location model (Section 1.3). Let

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

with \mathcal{F}_0 a subset of the collection of symmetric distributions (see (1.2)). Let $\hat{\mu}$ be an *equivariant estimator*, that is: the distribution of $\hat{\mu} - \mu$ does not depend on μ (see Chapter 9 for the formal definition of equivariance).

- 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.

which again does not depend on the parameter of interest.

- If $\mathcal{F}_0 := \{F_0 \text{ symmetric and continuous at } x=0\}$, 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 $p = 1/2$.

- Suppose now that X_1, \dots, X_n are the first n of an infinite sequence of i.i.d. random variables, and that

$$\mathcal{F}_0 := \{F_0 : \underbrace{\int x dF_0(x) = 0}_{\substack{\text{expectation} \\ \Rightarrow 0}}, \underbrace{\int x^2 dF_0(x) < \infty}_{\substack{\text{finite second moment}}} \}.$$

Variance in this case as the expectation $= \infty$.

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 $G = \Phi$.

6.3 Equivalence confidence sets and tests

Definition 6.3.1 A subset $I = I(\mathbf{X}) \subset \Gamma$, depending (only) on the data $\mathbf{X} = (X_1, \dots, X_n)$, is called a confidence set (Vertrauensbereich) for γ , at level $1-\alpha$, if

$$\mathbb{P}_\theta(\gamma \in I) \geq 1 - \alpha, \quad \forall \theta \in \Theta.$$

A confidence interval is of the form

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

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

Let for each $\gamma_0 \in \mathbb{R}$, $\phi(\mathbf{X}, \gamma_0) \in \{0, 1\}$ be a test at level α for the hypothesis $H_{\gamma_0} : \gamma = \gamma_0$.

Thus, we reject H_{γ_0} if and only if $\phi(\mathbf{X}, \gamma_0) = 1$, and

$$\mathbb{P}_{\theta: \gamma=\gamma_0}(\phi(\mathbf{X}, \gamma_0) = 1) \leq \alpha.$$

Then

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

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

Given the definition of the Null-hypothesis testing via pivot the confidence interval includes all of the parameters for which the test is rejected.

duality

62

CHAPTER 6. TESTS AND CONFIDENCE INTERVALS

Conversely, if $I(\mathbf{X})$ is a $(1 - \alpha)$ -confidence set for γ , then, for all γ_0 , the test $\phi(\mathbf{X}, \gamma_0)$ defined as

$$\phi(\mathbf{X}, \gamma_0) = \begin{cases} 1 & \text{if } \gamma_0 \notin I(\mathbf{X}) \\ 0 & \text{else} \end{cases}$$

is a test at level α of H_{γ_0} .

6.4 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. Recall that the power is of a test $\phi(X, \gamma_0)$ at a value θ with $g(\theta) \neq \gamma_0$ is $P_\theta(\phi(\mathbf{X}, \gamma_0) = 1)$.

power = probability of being out of the acceptance region when the parameter actually differs from the parameter.

6.5 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	(3, 4)	19

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.

6.5.1 Student's test

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

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

and the pooled sample variance

$$S^2 := \frac{1}{m+n-2} \left\{ \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (Y_j - \bar{Y})^2 \right\}. \quad \left. \begin{array}{l} \text{we assumed the} \\ \text{data to be} \\ \text{i.i.d.} \end{array} \right\}$$

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

$$\frac{\sigma^2}{n} + \frac{\sigma^2}{m} = \sigma^2 \left(\frac{n+m}{nm} \right). \quad \left. \begin{array}{l} \text{No Cov} \\ \text{by independence} \\ \text{assumption.} \end{array} \right\}$$

The normality assumption implies that

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

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 $\boxed{G \text{ is Student}(n+m-2)}$ (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 $\alpha = 0.05$.

6.5.2 Wilcoxon's 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

$$T = T^{\text{Wilcoxon}} := \sum_{i=1}^n R_i. \quad \text{Y} \quad \text{sum of ranks for the } X_i$$

Lemma:

$$E_{H_0} Z = \frac{n(N+1)}{2}$$

Proof:

$$P(R_i = k) = \frac{1}{N}$$

$$E_{H_0} R_i = \sum_{k=1}^N k \cdot \frac{1}{N} = \frac{N+1}{2}$$

$$E_{H_0} Z = \frac{n(N+1)}{2}$$

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

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

assuming that all X are smaller than Y .

correction for the rank.

so sum $\{1, \dots, n\}$

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

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

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

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

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

(see Theorem 6.5.1 below), and hence

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

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 6.5.1 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 4.1.4.

For two random variables X and Y , use the notation

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

when X and Y have the same distribution.

Theorem 6.5.1 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 $\mathbf{r} := (r_1, \dots, r_N)$,

not ordered

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

rather intuitive part.

As the distribution are equal the rank is totally random.

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

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

rearranging stuff

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

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

by def as this
is actually equal

by def
just

complex
notation

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

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

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

$$\mathbb{P}\left((Z_{(1)}, \dots, Z_{(n)}) \in A\right) = \underbrace{N! \mathbb{P}\left((Z_{q_1}, \dots, Z_{q_N}) \in A, Z_{q_1} < \dots < Z_{q_N}\right)}_{\text{as ordered now}} = N! \mathbb{P}\left((Z_{(1)}, \dots, Z_{(N)}) \in A, \mathbf{R} = \mathbf{r}\right),$$

independent of
permutation
 r, q will
always order.
so equality
holds for
all permutation

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

$$\mathbb{P}\left((Z_{(1)}, \dots, Z_{(N)}) \in A, \mathbf{R} = \mathbf{r}\right) = \frac{1}{N!} \mathbb{P}\left((Z_{(1)}, \dots, Z_{(n)}) \in A\right). \quad (6.1)$$

Take $A = \mathbb{R}^N$ to find that (6.1) implies

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

Plug this back into (6.1) to see that we have the product structure

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

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

¹Here is an example, with $N = 3$:

would
do!
The original

$$\begin{aligned} (z_1, z_2, z_3) &= (5, 6, 4) \\ (r_1, r_2, r_3) &= (2, 3, 1) \\ (q_1, q_2, q_3) &= (3, 1, 2) \end{aligned}$$

$$z_{q_1} = z_{(1)}$$

$$z_{q_2} = z_{(2)}$$

$$z_3 = z_{(3)}$$

at position
nearly
specified.

same as here

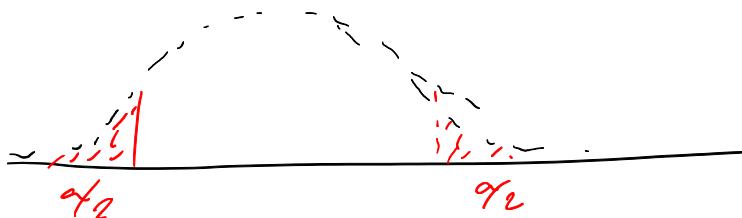
6.5.3 Comparison of Student's test and Wilcoxon's test

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 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.

Will appear in section 11.6.

Given the Wilcoxon test it is then possible to prove it by:

- (i) Either simulate from the null-distribution of Z .



Reject if $Z \in \underline{\underline{E}}$

- (ii) Calculate

$\text{Var}_{H_0}(Z)$ so that

$$Z_{\text{Standard-Norm}} = \frac{z - E_{H_0}(z)}{\sqrt{\text{Var}_{H_0}(z)}} \xrightarrow{D_{H_0}} N(0,1)$$

Notice despite that obs. dependent
it is possible to prove
that CLT holds

Check at the following Notes:

<https://web.stanford.edu/~lmackey/stats300a/doc/stats300a-fall15-lecture13.pdf>

Much more exhaustive than
the current.

Chapter 7

The Neyman Pearson Lemma and UMP tests

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 $X \in \mathcal{X}$, with distribution $P \in \mathcal{P}$, we consider the general testing problem for

$$\begin{aligned} H_0 : \theta &\in \Theta_0, \\ \text{against} \\ H_1 : \theta &\in \Theta_1. \end{aligned} \quad !$$

A (possibly randomized) test is some function $\phi : \mathcal{X} \rightarrow [0, 1]$. Fix some $\alpha \in [0, 1]$. We say that ϕ is a test at level α if

$$\sup_{\theta \in \Theta_0} E_\theta \phi(X) \leq \alpha. \quad \left\{ \begin{array}{l} = P_\theta(\phi=1) \\ \text{Type I error.} \end{array} \right.$$

Definition 7.0.1 A test ϕ is called Uniformly Most Powerful (UMP, (German: gleichmäßig mächtigst)) if

- ϕ has level α ,
- for all tests ϕ' with level α , it holds that $E_\theta \phi'(X) \leq E_\theta \phi(X) \forall \theta \in \Theta_1$.

7.1 The Neyman Pearson Lemma

We consider testing

$$H_0 : \theta = \theta_0$$

against the alternative

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

Define the risk $R(\theta, \phi)$ of a test ϕ as the probability of error of first and second kind:

$$R(\theta, \phi) := \begin{cases} E_\theta \phi(X), & \theta = \theta_0 \\ 1 - E_\theta \phi(X), & \theta = \theta_1 \end{cases} \quad \left\{ \begin{array}{l} \text{recall } \sum_{\theta} \int_{\mathcal{X}} \phi(x) dP_{\theta}(x) = 1 \\ \text{if } \theta = \theta_0, \text{ then } \phi \text{ rejects } H_0 \text{ iff } \phi(X) > \alpha \\ \text{if } \theta = \theta_1, \text{ then } \phi \text{ rejects } H_0 \text{ iff } \phi(X) < 1 - \alpha \end{array} \right.$$

by definition

the expectation 69

$E_\theta \phi(X, \theta_1) = 1 - E_\theta \phi(X, \theta_0)$

We want the test to reject H_0 when this is true the highest possible times.
i.e. we are looking for the highest power test.

power of any other test is lower.

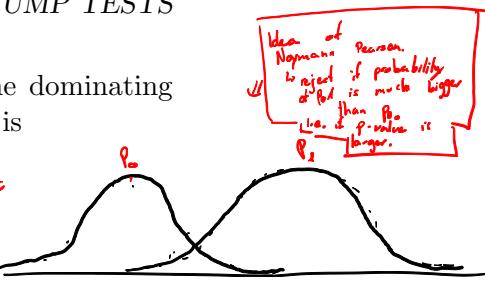
rejecting under alternative.

less power

less Type II error.

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

$$\phi_{NP} := \begin{cases} 1 & \text{if } p_1/p_0 > c \\ q & \text{if } p_1/p_0 = c \\ 0 & \text{if } p_1/p_0 < c \end{cases} \quad \begin{matrix} \rightarrow \text{reject if } n > pc \\ \text{or} \\ \text{if } p_1/p_0 > c \end{matrix}$$



Here $0 \leq q \leq 1$, and $0 \leq c < \infty$ are given constants.

Lemma 7.1.1 Neyman Pearson Lemma Let ϕ be some test. We have

$$\begin{aligned} R(\theta_1, \phi_{NP}) - R(\theta_1, \phi) &\leq c[R(\theta_0, \phi) - R(\theta_0, \phi_{NP})]. \end{aligned}$$

Difference in Power! Type II error actually

$R \rightarrow$ risk / 1 - Power. | decompose

This says Neyman Pearson is Minimum Power

Notice that it is not UMP yet as one needs to prove that the minimum power holds for all $\theta \neq \theta_0$.

such that the condition holds.

Proof.

$$\begin{aligned} R(\theta_1, \phi_{NP}) - R(\theta_1, \phi) &= \int (\phi - \phi_{NP}) p_1 \\ &= \int_{p_1/p_0 > c} (\phi - \phi_{NP}) p_1 + \int_{p_1/p_0 = c} (\phi - \phi_{NP}) p_1 + \int_{p_1/p_0 < c} (\phi - \phi_{NP}) p_1 \\ &\stackrel{\text{Def eq 7.1}}{\leq} c \int_{p_1/p_0 > c} (\phi - \phi_{NP}) p_0 + c \int_{p_1/p_0 = c} (\phi - \phi_{NP}) p_0 + c \int_{p_1/p_0 < c} (\phi - \phi_{NP}) p_0 \\ &\stackrel{\text{why?}}{\leq} c [R(\theta_0, \phi) - R(\theta_0, \phi_{NP})]. \end{aligned}$$

replace p_1 with $c \cdot p_0$

because of these characteristics and the interplay with ϕ ?

7.2 Uniformly most powerful tests

7.2.1 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

$$H_0 : \theta \in \left[\frac{1}{2}, 1 \right] =: \theta_0,$$

against the alternative

$$H_1 : \theta \in \left(\frac{1}{4}, \frac{1}{2} \right) =: \theta_1.$$

Important NP-Lemma

Note that

if you choose same level for the tests,

then by the Lemma above, flat type I error higher for other test

so some higher for NP-test

idea

for $\frac{P_{\theta_1}}{P_{\theta_0}} > c \quad \left\{ \begin{array}{l} \text{you reject if } \frac{P_{\theta_1}}{P_{\theta_0}} \text{ of successes too small for the given sat.} \\ \text{so you do it based on alternative.} \end{array} \right.$

71

7.2. UNIFORMLY MOST POWERFUL TESTS

does not depend on θ

i.e. $T(x)$ is sufficient to describe ratio P_{θ_1} and P_{θ_0}

Let $T := \sum_{i=1}^n X_i$ be the number of successes (T is a sufficient statistic), and consider the randomized test

$$\phi(T) := \begin{cases} 1 & \text{if } T < t_0 \\ q & \text{if } T = t_0, \\ 0 & \text{if } T > t_0 \end{cases}$$

why? As $T(x)$ is sufficient $P_{\theta_0}(x|T(x))$ does not depend on θ .

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 :

Thus, we take t_0 in such a way that

(i.e., $t_0 - 1 = q_{\inf}^G(\alpha)$) with q_{\inf}^G the quantile function defined in Section 6.1 and G the distribution function of T) and

$$\Rightarrow q = \frac{\alpha - P_{\theta_0}(T \leq t_0 - 1)}{P_{\theta_0}(T = t_0)}.$$

see big red above

Because $\phi = \phi_{NP}$ is the Neyman Pearson test, it is the most powerful test (at level α) (see the Neyman Pearson Lemma in Section 7.1). The power of the test is $\beta(\theta_1)$, where

$$\beta(\theta) := E_{\theta} \phi(T).$$

Numerical Example

Notice how interestingly, the sufficient statistic replaced the observed data.

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

-power = 1-TypeII

Notice now that $E_{\theta_1} \phi_{NP}(x) = P_{\theta_0}(x \leq 1) + q P_{\theta_1}(x=1)$
 $= \alpha$ iff q chosen like that
 solving above for q .

$$\beta(\theta_1) = P_{\theta_1}(T = 0) + q P_{\theta_1}(T = 1) \\ = \left(\frac{3}{4}\right)^7 + \frac{422}{546} \binom{7}{1} \left(\frac{3}{4}\right)^6 \left(\frac{1}{4}\right)$$

$$= 0.1335 + \frac{422}{546} 0.3114. \\ 1 - \left(P_{\theta_0}(T=0) + q \cdot P(T=1) \right)$$

gain in power due to randomization at q .

Problem 2

Consider now testing

$$H_0 : \theta_0 = \frac{1}{2},$$

against

$$H_1 : \theta < \frac{1}{2}.$$

*Have in
the entire
point. We
can ignore part 2
under Prob. 1
that we have
a pivot that
does not depend
on θ_0 or P_{θ_0} .
This was constructed
through θ_0 or P_{θ_0} .
So we
cannot ignore
Prob. 3.*

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 for the alternative $H_1 : \theta < \theta_0$.

Problem 3

We now want to test

$$H_0 : \theta \geq \frac{1}{2},$$

against the alternative

$$H_1 : \theta < \frac{1}{2}.$$

Recall the function

$$\beta(\theta) := E_\theta \phi(T).$$

The level of ϕ is defined as

$$\sup_{\theta \geq \frac{1}{2}} \beta(\theta).$$

*two components.
most powerful and
uniformity*

We have

$$\begin{aligned} \beta(\theta) &= P_\theta(T \leq t_0 - 1) + qP_\theta(T = t_0) \\ &= (1 - q)P_\theta(T \leq t_0 - 1) + qP_\theta(T \leq t_0). \end{aligned}$$

*complements case
+ red case
as now included
in here*

Observe that if $\theta_1 < \theta_0$, small values of T are more likely under P_{θ_1} than under P_{θ_0} :

$$P_{\theta_1}(T \leq t) > P_{\theta_0}(T \leq t), \forall t \in \{0, 1, \dots, n\}.$$

Thus, $\beta(\theta)$ is a decreasing function of θ . It follows that the level of ϕ is

$$\sup_{\theta \geq \frac{1}{2}} \beta(\theta) = \beta\left(\frac{1}{2}\right) = \alpha.$$

*highest likelihood i.e.
 $E_{\theta_0}(\phi)$*

*so if Ho changes
you have to reperform
the test.*

Hence, ϕ is uniformly most powerful for $H_0 : \theta \geq \frac{1}{2}$ against $H_1 : \theta < \frac{1}{2}$.

7.3 UMP tests and exponential families

We now study the situation where Θ is an interval in \mathbb{R} , and the testing problem is

$$H_0 : \theta \leq \theta_0,$$

against

$$H_1 : \theta > \theta_0.$$

We suppose that \mathcal{P} is dominated by a σ -finite measure ν .

Theorem 7.3.1 Suppose that \mathcal{P} is a one-dimensional exponential family

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

Assume moreover that $c(\theta)$ is a strictly increasing function of θ . Then a 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

$$\phi_{NP}(x) = \begin{cases} 1 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) > c_0 \\ q_0 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) = c_0 \\ 0 & \text{if } p_{\theta_1}(x)/p_{\theta_0}(x) < c_0 \end{cases},$$

where q_0 and c_0 are chosen in such a way that $E_{\theta_0}\phi_{NP}(X) = \alpha$. 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

$$\frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} = c \Leftrightarrow T(x) = t, \quad \begin{array}{l} \text{as given} \\ \text{the rest} \\ \text{determines} \\ \text{P}_{\theta_0}(x) \end{array}$$

where t is some constant (depending on c , θ_0 and θ_1). Therefore, $\phi = \phi_{NP}$. 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 $H_1 : \theta > \theta_0$.

You see
similar
reasoning
must
hold for
the previous
section.

We will now prove that $\beta(\theta) := E_\theta\phi(T)$ is increasing in θ . Let

to prove: $p_\theta(t) = \exp[c(\theta)t - d(\theta)]$

to show
Uniformity !

be the density of T with respect to dominating measure $\bar{\nu}$. For $\vartheta > \theta$

$$\frac{\bar{p}_\vartheta(t)}{\bar{p}_\theta(t)} = \exp \left[(c(\vartheta) - c(\theta))t - (d(\vartheta) - d(\theta)) \right],$$

$$\log \left(\frac{f_{\vartheta}}{f_\theta} \right) = 0 = \log \exp(S) = S(c(\vartheta) - c(\theta)) + -(d(\vartheta) - d(\theta))$$

which is increasing in t . Moreover, we have

$$(1) \quad \int \bar{p}_\vartheta d\bar{\nu} = \int \bar{p}_\theta d\bar{\nu} = 1. \quad \begin{array}{l} \text{because} \\ \text{or} \\ \text{dominating.} \end{array}$$

Therefore, there must be a point s_0 where the two densities cross:

$$(2) \quad \begin{cases} \frac{\bar{p}_\vartheta(t)}{\bar{p}_\theta(t)} \leq 1 & \text{for } t \leq s_0 \\ \frac{\bar{p}_\vartheta(t)}{\bar{p}_\theta(t)} \geq 1 & \text{for } t \geq s_0 \end{cases}.$$

so that
this holds.
where $\hat{x} < s_0$
 $\hat{x} > s_0$

$$\begin{aligned} & \left(c(\vartheta) - c(\theta) \right) \hat{x} - d(\vartheta) \neq 0 \\ & \left(c(\vartheta) - c(\theta) \right) \hat{x} - d(\vartheta) \neq 0 \end{aligned}$$

What happens if $c(\theta)$ decreases?

Then (1) holds,
(2) does not.

Opposite.

But then

$$\begin{aligned}\beta(\vartheta) - \beta(\theta) &= \int \phi(t)[\bar{p}_\vartheta(t) - \bar{p}_\theta(t)]d\bar{\nu}(t) \\ \text{Therefore} &= \int_{t \leq s_0} \phi(t)[\bar{p}_\vartheta(t) - \bar{p}_\theta(t)]d\bar{\nu}(t) + \int_{t \geq s_0} \phi(t)[\bar{p}_\vartheta(t) - \bar{p}_\theta(t)]d\bar{\nu}(t) \\ &\geq \phi(s_0) \int [\bar{p}_\vartheta(t) - \bar{p}_\theta(t)]d\bar{\nu}(t) = 0.\end{aligned}$$

min positive

So indeed $\beta(\theta)$ is increasing in θ . \Rightarrow as $\alpha > 0$

But then

$$\sup_{\theta \leq \theta_0} \beta(\theta) = \boxed{\beta(\theta_0)} = \alpha.$$

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. \square

Example 7.3.1 Test for the variance of the normal distribution

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

normal dist in exponential form

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

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 $T(\mathbf{X})$. Note that under H_0 , the statistic $T(\mathbf{X})/\sigma_0^2$ has a χ^2 -distribution with n degrees of freedom, the χ^2 -distribution (see Section 12.2 for a definition). So we can find the critical value from the quantile of the χ^2 -distribution.

from previous section

7.4 One- and two-sided tests: an example with the Bernoulli distribution

Let X_1, \dots, X_n be an i.i.d. sample from the $\text{Bernoulli}(\theta)$ -distribution, $0 < \theta < 1$.

Then

$$p_\theta(x_1, \dots, x_n) = \exp \left[\log \left(\frac{\theta}{1-\theta} \right) \sum_{i=1}^n x_i + n \log(1-\theta) \right].$$

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 .

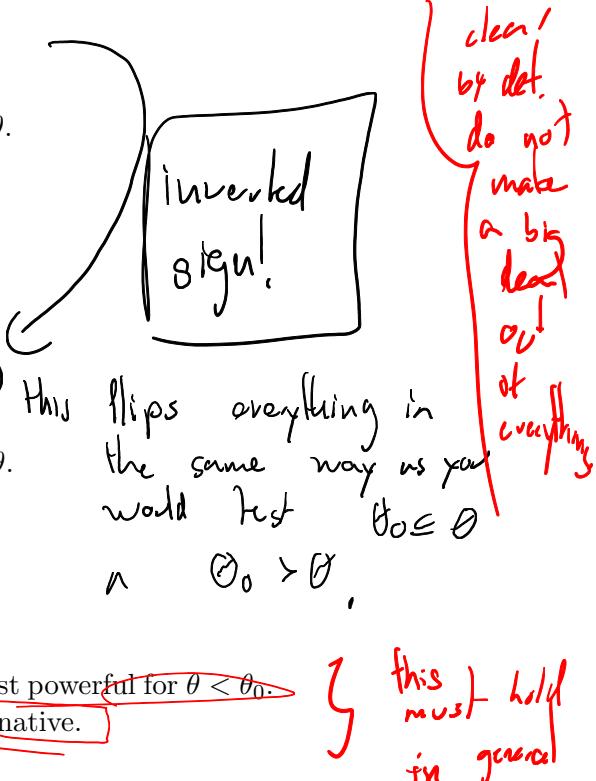
Two-sided alternative

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

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.



7.5 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 $\alpha (< 1)$.

As we have seen in Section 7.4, 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.

1 → reject

	reject	accept	
Ho true	type I	power	better to take opposite test.
Ho false	Power	Type II	if this does not hold for you

Rejections when wrong



CHAPTER 7. THE NEYMAN PEARSON LEMMA AND UMP TESTS

Definition 7.5.1 A test ϕ is called unbiased (German unverfälscht) if for all $\theta \in \Theta_0$ and all $\vartheta \in \Theta_1$,

Type I Power

$$E_\theta \phi(X) \leq E_\vartheta \phi(X).$$

Power of rejecting under the null is smaller than the power of rejecting under the alternative.

Definition 7.5.2 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 $H_1 : \theta \neq \theta_0$.

desired property; we want to reject if $\theta \in \Theta_1$.

The following theorem presents the UMPU test. We omit the proof (see e.g. Lehmann (1986)).

Theorem 7.5.1 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 a 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

$$E_{\theta_0} \phi(X) = \alpha, \quad \left. \frac{d}{d\theta} E_\theta \phi(X) \right|_{\theta=\theta_0} = 0. \quad \begin{array}{l} \text{Unbiased condition.} \\ \text{you are at minimum risk for } \theta = \theta_0. \end{array}$$

As if depends on you are which side testing; i.e. $H_0: \theta \geq \theta_0$ or $H_0: \theta \leq \theta_0$.

Note Let ϕ_R a right-sided test as defined Theorem 7.3.1 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 7.5.1 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

$$\beta(\theta) = \beta_R(\theta) + \beta_L(\theta).$$

so the sum of the two one-sided risks should be minimized.

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 7.3.1, for a test $\tilde{\phi}$ depending only on the sufficient statistic T ,

$$E_\theta \tilde{\phi}(T) = \int \tilde{\phi}(t) \exp[c(\theta)t - d(\theta)] d\bar{\nu}(t).$$

↑ see bernoulli example it not clear.

if a Null hypothesis follows, imply that $E_\theta \phi(T) \geq E_{\theta_0} \phi(X)$

usual slot then it follows, imply that $E_\theta \phi(T) \geq E_{\theta_0} \phi(X)$

level off

level off

$$\frac{d}{d\theta} E(\tilde{\phi}(T)) = E(\tilde{\phi}'(T))$$

Recall

Hence, assuming we can take the differentiation inside the integral,

$$\begin{aligned} \frac{d}{d\theta} E_\theta \tilde{\phi}(T) &= \int \tilde{\phi}(t) \exp[c(\theta)t - d(\theta)] (\dot{c}(\theta)t - \dot{d}(\theta)) d\nu(t) \\ &= \dot{c}(\theta) \text{cov}_\theta(\tilde{\phi}(T), T). \end{aligned}$$

?

The UMPU test sets this to zero. We leave the interpretation to the reader...

Example 7.5.1 Two-sided test for the mean of the normal distribution

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,

$$\frac{d}{d\mu} E_\mu \phi(T) = \frac{n}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) - \frac{n}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right),$$

So putting

$$\frac{d}{d\mu} E_\mu \phi(T) \Big|_{\mu=\mu_0} = 0,$$

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

inver^t
this:

$$t_R - n\mu_0 = \sqrt{n}\sigma_0\Phi^{-1}(1 - \alpha/2), \quad t_L - n\mu_0 = -\sqrt{n}\sigma_0\Phi^{-1}(1 - \alpha/2).$$

7.6 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 7.6.1 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.



Sketch of proof.

Let $\bar{p}_\theta(t_1, t_2)$ be the density of (T_1, T_2) with respect to dominating measure $\bar{\nu}$:

$$\bar{p}_\theta(t_1, t_2) := \exp[\beta t_1 + \gamma t_2 - d(\theta)]\bar{h}(t_1, t_2).$$

We assume $\bar{\nu}(t_1, t_2) = \bar{\nu}_1(t_1)\bar{\nu}_2(t_2)$ is a product measure. The conditional density of T_2 given $T_1 = t_1$ is then

$$\begin{aligned} \bar{p}_\theta(t_2|t_1) &= \frac{\exp[\beta t_1 + \gamma t_2 - d(\theta)]\bar{h}(t_1, t_2)}{\int_{s_2} \exp[\beta t_1 + \gamma s_2 - d(\theta)]\bar{h}(t_1, s_2)d\bar{\nu}_2(s_2)} \\ &= \exp[\gamma t_2 - d(\gamma|t_1)]\bar{h}(t_1, t_2), \end{aligned}$$

where

$$d(\gamma|t_1) := \log \left(\int_{s_2} \exp[\gamma s_2]\bar{h}(t_1, s_2)d\bar{\nu}_2(s_2) \right).$$

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,

$$\sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi(T) = \sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \underbrace{E_{\gamma}(\phi(T)|T_1)}_{\leq \alpha} \leq \alpha.$$

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*

$$\alpha' := \sup_{\beta} \sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi'(T) \leq \alpha,$$

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

$$\alpha' = \sup_{\beta} \sup_{\gamma \leq \gamma_0} E_{(\beta, \gamma)} \phi'(T)$$

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

$$E_\gamma(\phi'(T)|T_1) \leq E_\gamma(\phi(T)|T_1), \forall \gamma > \gamma_0.$$

Thus also

$$E_{(\beta, \gamma)} \phi'(T) \leq E_{(\beta, \gamma)} \phi(T), \forall \beta, \gamma > \gamma_0.$$

□

Example 7.6.1 Comparing the means of two Poissons

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

$$\beta := \log(\mu), \gamma := \log(\lambda/\mu).$$

The testing problem is equivalent to

$H_0 : \gamma \leq \gamma_0$,

against the alternative

$H_1 : \gamma > \gamma_0$,

where $\gamma_0 := 0$.

The density is

$$\mathbf{p}_\theta(x_1, \dots, x_n, y_1, \dots, y_m)$$

$$\begin{aligned}
&= \exp \left[\log(\lambda) \sum_{i=1}^n x_i + \log(\mu) \sum_{j=1}^m y_j - n\lambda - m\mu \right] \prod_{i=1}^n \frac{1}{x_i!} \prod_{j=1}^m \frac{1}{y_j!} \\
&= \exp \left[\log(\mu) \left(\sum_{i=1}^n x_i + \sum_{j=1}^m y_j \right) + \log(\lambda/\mu) \sum_{i=1}^n x_i - n\lambda - m\mu \right] h(\mathbf{x}, \mathbf{y}) \\
&= \exp[\beta T_1(\mathbf{x}, \mathbf{y}) + \gamma T_2(\mathbf{x}) - d(\theta)] h(\mathbf{x}, \mathbf{y}),
\end{aligned}$$

where

$$T_1(\mathbf{X}, \mathbf{Y}) := \sum_{i=1}^n X_i + \sum_{j=1}^m Y_j, \quad T_2(\mathbf{X}) := \sum_{i=1}^n X_i,$$

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 8

Comparison of estimators

One can compare estimators in terms of their *risk*. This is done in Sections 8.1, 8.2 and 8.3. Section 8.4 briefly addresses sensitivity and robustness and Section 8.5 discusses computational aspects. These last two sections do not present the details.

8.1 Definition of risk

Consider a random variable X with distribution P_θ , $\theta \in \Theta$. Let $T = T(X)$ be an estimator of a parameter of interest $\gamma = g(\theta)$. A risk function $R(\cdot, \cdot)$ measures the loss due to the error of the estimator. The *risk* depends on the unknown parameter θ and on the estimator. We define the risk as

$$R(\theta, T) := \mathbb{E}_\theta(L(\theta, T(X)))$$

where $L(\cdot, \cdot)$ is a given so-called *loss function*¹. A more detailed description is given in Chapter 10.

Example 8.1.1 Risk of a test

Consider the testing problem

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

against

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

Let $\phi(X) \in [0, 1]$ be a test.

The risk of the test can then be defined as the probability of an error, i.e.

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

Example 8.1.2 Risk of an estimator

In the case $\gamma \in \mathbb{R}$ an important risk measure is the mean square error

$$R(\theta, T) := \mathbb{E}_\theta(T(X) - g(\theta))^2 =: \text{MSE}_\theta(T).$$

¹Note that the quantity $L(\theta, T(X))$ is random. Note also that in the notation of risk $R(\theta, T)$, the symbol T stands for the map T .

8.2 Risk and sufficiency

Let $S = S(X)$ be sufficient. Knowing the sufficient statistic S one can forget about the original data X without losing information. Indeed, the following lemma says that any decision based on the original data X can be replaced by a randomized one which depends only on S and which has the same risk.

Lemma 8.2.1 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

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

His is a random variable $\delta(S) = d(X^*)$, where X^* are generated data.

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

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

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

So this is basically the proof and the \square .
understanding about a sufficient statistic is given.

8.3 Rao-Blackwell

The Lemma of Rao-Blackwell says that in the case of convex loss an estimator based on the original data X can be replaced by one based only on S without increasing the risk. Randomization is not needed here.

Lemma 8.3.1 (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

This guarantees that the union of two convex functions is convex.
 This is why randomization is not needed here.

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

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

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

Hence, $\forall \theta$,

$$\begin{aligned} E\left(L\left(\theta, d(X)\right) \middle| S = s\right) &\geq L\left(\theta, E\left(d(X) | S = s\right)\right) \quad (1) \\ &= L(\theta, d'(s)). \end{aligned}$$

By the iterated expectations lemma, we arrive at

$$\begin{aligned}
 R(\theta, d) &= E_\theta L(\theta, d(X)) \\
 &= E_\theta E\left(L\left(\theta, d(X)\right) \middle| S\right) \quad \text{by} \\
 &\geq \underbrace{E_\theta L(\theta, d'(S))}_{R(\theta, d')} \quad \square
 \end{aligned}$$

Example 8.3.1 Mean square error

Let T be an estimator of $g(\theta) \in \mathbb{R}$ and let

$$R(\theta, T) := \mathbb{E}_\theta(T(X) - g(\theta))^2 =: \text{MSE}_\theta(T).$$

Let S be sufficient and $\tilde{T} := \mathbb{E}(T|S)$. Then by the Rao-Blackwell Lemma

$$\text{MSE}_\theta(T) \geq \text{MSE}_\theta(\tilde{T}) \iff R(\theta, \tilde{T}) < R(\theta, T) \forall \theta.$$

The mean square error can be decomposed in the variance term and the squared bias term. Since $E\tilde{T} = ET$ by the iterated expectations lemma, we thus have

→ relationship given that $E(\tilde{T}) = E(T)$ and therefore bias equality holds.

So that $\text{MSE}_\theta(T) \geq \text{MSE}_\theta(\tilde{T})$ since $\text{Var}_\theta(\tilde{T}) \leq \text{var}_\theta(T), \forall \theta$.

Compare with Lemma 5.2.2 and the result of Lehmann-Scheffé in Section 5.3.

8.4 Sensitivity and robustness

We can compare estimators with respect to their sensitivity to large errors in the data. Let X_1, \dots, X_n be i.i.d. copies of a random variable X . Let $T_n = T_n(X_1, \dots, X_n)$ be a real-valued estimator defined for each n , and symmetric in X_1, \dots, X_n .

Influence of a single additional observation

The influence function is

$$l(x) := T_{n+1}(X_1, \dots, X_n, x) - T_n(X_1, \dots, X_n), \quad x \in \mathbb{R}.$$

Break down point

Let for $m < n$.

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

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

$$\epsilon^* := \min\{m : \epsilon(m) = \infty\}/n. \quad \left. \begin{array}{l} \text{share of} \\ \text{outliers needed} \\ \text{in order for the} \\ \text{statistic to break down} \end{array} \right\}$$

An estimator is called robust if it has a bounded influence function and/or a large breakdown point.

8.5 Computational aspects

Today, the data are often high-dimensional and the number of parameters p is also very large. Maximum likelihood estimation for example requires maximization of a function of p variables and this can be very hard if p is large. The more so if the likelihood is not concave, or if there are e.g. some parameters are integer valued etc. We will moreover examine in Chapter 10 Bayesian theory. Then one needs to find so-called “posterior distributions”, which is typically computationally very hard (this is where MCMC (Monte Carlo Markov Chain) algorithms come in). Clearly, an estimator which cannot be computed (say in polynomial time) is of little practical value.

Chapter 9

Equivariant statistics

As we have seen in Chapter 5 for instance, it can be useful to restrict attention to a collection of statistics satisfying certain desirable properties. In Chapter 5, 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

$$X_i = \theta + \epsilon_i, \quad i = 1, \dots, n. \quad \theta \text{ as location parameter}$$

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} . We assume $\epsilon_1, \dots, \epsilon_n$ are i.i.d. with a known density $p_0(\cdot)$.

Location-scale model

Here $\theta = (\mu, \sigma)$, with $\mu \in \mathbb{R}$ a location parameter and $\sigma > 0$ a scale parameter. We assume

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

The parameter space Θ and action space \mathcal{A} are both $\mathbb{R} \times (0, \infty)$. We assume $\epsilon_1, \dots, \epsilon_n$ are i.i.d. with a known density $p_0(\cdot)$.

9.1 Equivariance in the location model

Definition 9.1.1 A statistic $T = T(\mathbf{X})$ is called location equivariant if for all constants $c \in \mathbb{R}$ and all $\mathbf{x} = (x_1, \dots, x_n)$,

$$T(x_1 + c, \dots, x_n + c) = T(x_1, \dots, x_n) + c.$$

*y transforming $p(\cdot)$
before or after makes no difference*

Examples

$$T = \begin{cases} \bar{X} & (n \text{ odd}) \\ X_{(\frac{n+1}{2})} & (n \text{ even}) \\ \dots \end{cases}$$

losses based on variance

Definition 9.1.2 A loss function $L(\theta, a)$ is called location invariant if for all $c \in \mathbb{R}$,

$$L(\theta + c, a + c) = L(\theta, a), (\theta, a) \in \mathbb{R}^2.$$

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 9.1.1 If T is equivariant (and $L(\theta, a)$ is invariant), then ! **Important**

$$\begin{aligned} R(\theta, T) &= E_\theta L(\theta, T(\mathbf{X})) = E_\theta L(0, T(\mathbf{X}) - \theta) \\ &= E_\theta L(0, T(\mathbf{X} - \theta)) = E_\theta L_0[T(\varepsilon)], \end{aligned}$$

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

$$R(\theta, T) = E_\theta L_0[T(\mathbf{X})] = R(0, T). \quad \left\{ \text{this generally true.} \right.$$

Definition 9.1.3 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,

$$R(0, T) = \min_{d \text{ equivariant}} R(0, d).$$

so the risk function with lower risk where d is equivariant

9.1.1 Construction of the UMRE estimator

Lemma 9.1.1 Let $Y_i := X_i - X_n$, $i = 1, \dots, n$, and $\mathbf{Y} := (Y_1, \dots, Y_n)$. We have

$$T \text{ equivariant} \Leftrightarrow T(\mathbf{X}) = T(\mathbf{Y}) + X_n.$$

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

this proves equivariance

given equivariance
so that $= T(x) - x_n + y_n$
 $= T(x) - x_n + y_n$
as you would add
also to the subtraction.

Theorem 9.1.1 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 \mathbf{Y} and its distribution 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

$$\begin{aligned} \text{Hence } T(\mathbf{X}) - \theta &= T(\mathbf{Y}) + \epsilon_n. \\ &\quad \xrightarrow{\text{as } T(\mathbf{Y}) \text{ is } \text{equivariant}} \text{by the definition of } \mathbf{X} \text{ in the location model. Here this was replaced.} \\ R(0, T) &= E L_0(T(\mathbf{Y}) + \epsilon_n) = E E \left[L_0(T(\mathbf{Y}) + \epsilon_n) \middle| \mathbf{Y} \right]. \\ \text{But } R(0, T) &= N(0, T(\mathbf{Y})) \quad \xrightarrow{\substack{\text{law of iterated expectation} \\ \text{minimum}}} \\ E \left[L_0(T(\mathbf{Y}) + \epsilon_n) \middle| \mathbf{Y} \right] &\geq \min_v E \left[L_0(v + \epsilon_n) \middle| \mathbf{Y} \right] \quad \{ \text{by def of the minimum function.} \} \\ &= E \left[L_0(T^*(\mathbf{Y}) + \epsilon_n) \middle| \mathbf{Y} \right]. \quad \{ \text{by def in box above} \} \end{aligned}$$

Hence,

$$R(0, T) \geq E E \left[L_0(T^*(\mathbf{Y}) + \epsilon_n) \middle| \mathbf{Y} \right] = R(0, T^*). \quad \{ \text{and proves that } T^* \text{ is UMRE by def of UMRE.} \}$$

9.1.2 Quadratic loss: the Pitman estimator

Corollary 9.1.2 If we take quadratic loss

Recall that by the previous proof in this section and use of invariant loss and equivariant T , you can always write $L(\theta, a) := (a - \theta)^2$, we get $L_0[a] = a^2$, and so, for $\mathbf{Y} = \mathbf{X} - X_n$,

Recall that $L(\theta, a) = (a - \theta)^2$

and hence replacing θ by v and a by $v + \epsilon_n$ at it.

$$T^*(\mathbf{Y}) = \arg \min_v E \left[(v + \epsilon_n)^2 \middle| \mathbf{Y} \right]$$

$$= -E(\epsilon_n | \mathbf{Y}),$$

$$T^*(\mathbf{X}) = X_n - E(\epsilon_n | \mathbf{Y}).$$

This estimator is called the Pitman estimator.

$$\begin{aligned} E(a - \theta)^2 &= E(a^2 - 2a\theta + \theta^2) \\ -Ea &= E(a^2) - 2E(a)\theta + \theta^2 \\ \frac{\partial E}{\partial \theta} &= E(-2a) + 2E(a) = 0 \\ a &= E(a) \end{aligned}$$

Notice this is invariant as inserted L but same as before

this is UMRE

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

$$f_Y(y) = \int f(y+z, z) dz. \quad \int f(y+x-y, z) dz = \int f(x, z) dz$$

Lemma 9.1.2 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

$$\underline{T^*(\mathbf{X})} = \frac{\int z \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}{\int \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}.$$

Proof. Let $\mathbf{Y} = \mathbf{X} - X_n$. The random vector \mathbf{Y} has density

$$\underbrace{f_{\mathbf{Y}}(y_1, \dots, y_{n-1}, 0)}_{\text{by def}} = \int \mathbf{p}_0(y_1 + z, \dots, y_{n-1} + z, z) dz.$$

So the density of ϵ_n given $\mathbf{Y} = \mathbf{y} = (y_1, \dots, y_{n-1}, 0)$ is

$$f_{\epsilon_n}(u) = \frac{p_0(y_1 + u, \dots, y_{n-1} + u, u)}{\int p_0(y_1 + z, \dots, y_{n-1} + z, z) dz}.$$

It follows that

$$E(\epsilon_n | \mathbf{y}) = \frac{\int u p_0(y_1 + u, \dots, y_{n-1} + u, u) du}{\int p_0(y_1 + z, \dots, y_{n-1} + z, z) dz}.$$

Thus

$$\begin{aligned}
 E(\epsilon_n | \mathbf{Y}) &= \frac{\int u \mathbf{p}_0(Y_1 + u, \dots, Y_{n-1} + u, u) du}{\int \mathbf{p}_0(Y_1 + z, \dots, Y_{n-1} + z, z) dz} \quad \text{inverkd} \\
 &= \frac{\int u \mathbf{p}_0(X_1 - X_n + u, \dots, X_{n-1} - X_n + u, u) du}{\int \mathbf{p}_0(X_1 - X_n + z, \dots, X_{n-1} - X_n + z, z) dz} \\
 &= X_n \left(\frac{\int \mathbf{p}_0(X_1 - z, \dots, X_{n-1} - z, X_n - z) dz}{\int \mathbf{p}_0(X_1 - z, \dots, X_{n-1} - z, X_n - z) dz} \right).
 \end{aligned}$$

Finally, recall that $T^*(\mathbf{X}) = X_n - E(\overline{\epsilon_n | \mathbf{Y}})$.

Change of
variable here.
Inserting:
 $z = x_n - v$

Example 9.1.1 Uniform distribution with unknown midpoint

Suppose X_1, \dots, X_n are i.i.d. $\text{Uniform}[\theta - 1/2, \theta + 1/2]$, $\theta \in \mathbb{R}$. Then

We have

$$\max_{1 \leq i \leq n} |x_i - z| \leq 1/2 \iff \underline{x_{(n)}} - 1/2 \leq z \leq \overline{x_{(1)}} + 1/2.$$

S₀

$$p_0(x_1 - z, \dots, x_n - z) = 1\{x_{(n)} - 1/2 \leq z \leq x_{(1)} + 1/2\}.$$

Thus, writing

$$T_1 := X_{(n)} - 1/2, \quad T_2 := X_{(1)} + 1/2,$$

the UMBE estimator T^* is

$$T^* = \left(\int_{T_1}^{T_2} \mathbf{1} \otimes dz \right) \Big/ \left(\int_{T_1}^{T_2} \mathbf{1} \cdot dz \right) = \frac{T_1 + T_2}{2} = \frac{X_{(1)} + X_{(n)}}{2}.$$

as PPF indicator funct.

||| D Want
o this
form in
the end.

9.1.3 Invariant statistics

*not equivariant*We now consider more general invariant statistics \mathbf{Y} .**Definition 9.1.4** A map $\mathbf{Y} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ 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}' .)

The map is
reproducible
= constant
up to dependent
on \mathbf{x}, \mathbf{x}' .

Might be in
higher dim and
different for x_1, \dots, x_n

Notice that
when this is R^n ,
this is before
we do any
computation.

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:

$\mathbf{x} = \mathbf{x}' + (x_n - x'_n)$ and this \mathbf{x}'
 $d(\mathbf{x} + c) = d(\mathbf{x}) + c$ and this
then c

Example Let $d(\mathbf{X})$ be equivariant. Then $\mathbf{Y} := \mathbf{X} - d(\mathbf{X})$ is maximal invariant.**Theorem 9.1.2** Suppose that $d(\mathbf{X})$ is equivariant. Let $\mathbf{Y} := \mathbf{X} - d(\mathbf{X})$, and

$$T^*(\mathbf{Y}) := \arg \min_v E \left[L_0(v + d(\varepsilon)) \mid \mathbf{Y} \right].$$

Then

$$T^*(\mathbf{X}) := T^*(\mathbf{Y}) + d(\mathbf{X})$$

is UMRE.

$$\mathbf{x} - d(\mathbf{x})$$

Proof. Let T be an equivariant estimator. Then

$$\begin{aligned} T(\mathbf{X}) &= T(\mathbf{X} - d(\mathbf{X})) + d(\mathbf{X}) \\ &= T(\mathbf{Y}) + d(\mathbf{X}). \end{aligned}$$

Follows from here
where $\varepsilon \rightarrow X$.

Hence

$$\begin{aligned} E \left[L_0(T(\varepsilon)) \mid \mathbf{Y} \right] &= E \left[L_0(T(\mathbf{Y}) + d(\varepsilon)) \mid \mathbf{Y} \right] \\ &\geq \min_v E \left[L_0(v + d(\varepsilon)) \mid \mathbf{Y} \right]. \end{aligned}$$

Now, use the iterated expectation lemma.

so if estim T \square equivariant
 $T(\mathbf{x}) = E_0(T(\mathbf{x}) \mid \mathbf{x} - T(\mathbf{x}))$
 $= T(\mathbf{x} - E_0(T(\mathbf{x}) \mid \mathbf{x} - T(\mathbf{x}))$

9.1.4 Quadratic loss and Basu's Lemma

For quadratic loss ($L_0[a] = a^2$), the definition of $T^*(\mathbf{Y})$ in the above theorem is

$$\begin{aligned} T^*(\mathbf{Y}) &= -E(d(\varepsilon) \mid \mathbf{Y}) = -E_0(d(\mathbf{X}) \mid \mathbf{X} - d(\mathbf{X})), \\ \text{so that } & \text{Replacing } \mathbf{x} \text{ with } \mathbf{x} - d(\mathbf{x}) \Rightarrow T^*(\mathbf{X}) = d(\mathbf{X}) - E_0(d(\mathbf{X}) \mid \mathbf{X} - d(\mathbf{X})). \end{aligned}$$

Recall when $\theta = 0$
 $x = \varepsilon$

So for a equivariant estimator T , we have

$$T \text{ is UMRE} \Leftrightarrow E_0(T(\mathbf{X}) \mid \mathbf{X} - T(\mathbf{X})) = 0.$$

??

$$\begin{aligned} \text{clear } & T(\mathbf{x}) = T(\mathbf{y}) + d(\mathbf{x}) \\ \text{by equivariant } & d(\mathbf{x}) = l(\mathbf{x} - E_0(l(\mathbf{x}) \mid \mathbf{x} - d(\mathbf{x}))) \text{ with } d(\mathbf{x}) = T(\mathbf{x}) \quad T(\mathbf{y}) = 0 \end{aligned}$$

$d(\mathbf{x})$
is
equivariant.

Check
notes
part.

location note and you have your result.

$$\begin{aligned}
 \text{Reason } E_0 T^*(x) &\rightarrow \text{location note} \\
 &= E_0 T(\theta + \epsilon) \\
 &= E_0 (\theta + T(\epsilon)) \rightarrow \text{equivariance} \\
 &= \theta + E_0 T^*(\epsilon) = \theta + \theta. \quad \text{by law iterated expectation} \\
 &\Rightarrow E_0 T^*(\epsilon) = 0 \quad \text{so } E_0 T(x) = 0
 \end{aligned}$$

92

CHAPTER 9. EQUIVARIANT STATISTICS

From the right hand side, we conclude that $E_0 T = 0$ and hence $E_\theta(T) = \theta$. 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

$$E_0(T(\mathbf{X})|\mathbf{X} - T(\mathbf{X})) = E_0 T(\mathbf{X}) = 0.$$

So then T is UMRE \rightarrow as it was already equivariant and unbiased.

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

$$E_\theta h(T) = 0, \forall \theta, \rightarrow \text{due to tower rule}$$

we conclude from the completeness of T that

$$h(T) = 0, P_\theta\text{-a.s., } \forall \theta,$$

in other words,

$$P(Y \in A|T) = P(Y \in A), P_\theta\text{-a.s., } \forall \theta.$$

Since A was arbitrary, we thus have that the conditional distribution of Y given T is equal to the unconditional distribution:

$$P(Y \in \cdot|T) = P(Y \in \cdot), P_\theta\text{-a.s., } \forall \theta,$$

$$P(Y|T) = \frac{P(Y|T)}{P(T)}$$

that is, for all θ , T and Y are independent under P_θ . \square

Basu's lemma is intriguing: it proves a probabilistic property (independence) via statistical concepts.

Example 9.1.2 UMRE estimator for the mean of the normal distribution: σ^2 known

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. \rightarrow As then risk unbiased as previously shown.

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 $\mathbf{X} - \bar{X}$.

means these two must be equal.
and by proof must be indep.

the direction

Needs suff
↑ and complete
Basu's start
Y being independent
Lemma of

every straightforward

by exponential
Theorem

See exercise sheet

9.2 Equivariance in the location-scale model *

Location-scale model

We assume

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

The unknown parameter is $\theta = (\mu, \sigma)$, with $\mu \in \mathbb{R}$ a location parameter and $\sigma > 0$ a scale parameter. The parameter space Θ and action space \mathcal{A} are both $\mathbb{R} \times \mathbb{R}_+$ ($\mathbb{R}_+ := (0, \infty)$). The distribution of $\varepsilon = (\epsilon_1, \dots, \epsilon_n)$ is assumed to be known.

Definition 9.2.1 A statistic $T = T(\mathbf{X}) = (T_1(\mathbf{X}), T_2(\mathbf{X}))$ is called location-scale equivariant if for all constants $b \in \mathbb{R}$, $c \in \mathbb{R}_+$, and all $\mathbf{x} = (x_1, \dots, x_n)$,

$$T(b + cx_1, \dots, b + cx_n) = b + cT(x_1, \dots, x_n)$$

and

$$T_2(b + cx_1, \dots, b + cx_n) = cT_2(x_1, \dots, x_n).$$

Definition 9.2.2 A loss function $L(\mu, \sigma, a_1, a_2)$ is called location-scale invariant if for all $(\mu, a_1, b) \in \mathbb{R}^3$, $(\sigma, a_2, c) \in \mathbb{R}_+^3$

$$L(b + c\mu, c\sigma, b + ca_1, ca_2) = L(\mu, \sigma, a_1, a_2).$$

In this section we abbreviate location-scale equivariance (invariance) to simply equivariance (invariance), and we assume throughout that the loss $L(\theta, a)$ is invariant.

Corollary 9.2.1 If T is equivariant (and $L(\theta, a)$ is invariant), then

$$\begin{aligned} R(\theta, T) &= E_\theta L(\mu, \sigma, T_1(\mathbf{X}), T_2(\mathbf{X})) \\ &= E_\theta L\left(0, 1, \frac{T_1(\mathbf{X}) - \mu}{\sigma}, \frac{T_2(\mathbf{X})}{\sigma}\right) \\ &= E_\theta L\left(0, 1, T_1(\varepsilon), T_2(\varepsilon)\right) = E_\theta L_0(T(\varepsilon)), \end{aligned}$$

where $L_0(a_1, a_2) := L(0, 1, a_1, a_2)$. We conclude that the risk does not depend on θ . We may therefore omit the subscript θ in the last expression:

$$R(\theta, T) = EL_0(T(\varepsilon)).$$

Definition 9.2.3 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,

$$R(0, 1, T_1, T_2) = \min_{d \text{ equivariant}} R(0, 1, d_1, d_2).$$

9.2.1 Construction of the UMRE estimator *

Theorem 9.2.1 Suppose that $d(\mathbf{X})$ is equivariant. Let

$$\mathbf{Y} := \frac{\mathbf{X} - d_1(\mathbf{X})}{d_2(\mathbf{X})},$$

and

$$T^*(\mathbf{Y}) := \arg \min_{a_1 \in \mathbb{R}, a_2 \in \mathbb{R}_+} E \left[L_0 \left(d_1(\varepsilon) + d_2(\varepsilon)a_1, d_2(\varepsilon)a_2 \right) \middle| \mathbf{Y} \right].$$

Then

$$T^*(\mathbf{X}) := \begin{pmatrix} d_1(\mathbf{X}) + d_2(\mathbf{X})T_1^*(\mathbf{Y}) \\ d_2(\mathbf{X})T_2^*(\mathbf{Y}) \end{pmatrix}$$

is UMRE.

Proof. We have

$$\mathbf{Y} = \frac{\mathbf{X} - d_1(\mathbf{X})}{d_2(\mathbf{X})} = \frac{\varepsilon - d_1(\varepsilon)}{d_2(\varepsilon)}.$$

So

$$\varepsilon = d_1(\varepsilon) + d_2(\varepsilon)\mathbf{Y}.$$

Let T be an equivariant estimator. Then

$$\begin{aligned} & EL_0 \left(T_1(\varepsilon), T_2(\varepsilon) \right) \\ &= EL_0 \left(T_1(d_1(\varepsilon) + d_2(\varepsilon)\mathbf{Y}), T_2(d_1(\varepsilon) + d_2(\varepsilon)\mathbf{Y}) \right) \\ &= EL_0 \left(d_1(\varepsilon) + d_2(\varepsilon)T_1(\mathbf{Y}), d_2(\varepsilon)T_2(\mathbf{Y}) \right) \\ &= EE \left[L_0 \left(d_1(\varepsilon) + d_2(\varepsilon)T_1(\mathbf{Y}), d_2(\varepsilon)T_2(\mathbf{Y}) \right) \middle| \mathbf{Y} \right] \\ &\geq E \min_{a_1 \in \mathbb{R}, a_2 \in \mathbb{R}_+} E \left[L_0 \left(d_1(\varepsilon) + d_2(\varepsilon)a_1, d_2(\varepsilon)a_2 \right) \middle| \mathbf{Y} \right] \\ &= EE \left[L_0 \left(d_1(\varepsilon) + d_2(\varepsilon)T_1^*(\mathbf{Y}), d_2(\varepsilon)T_2^*(\mathbf{Y}) \right) \middle| \mathbf{Y} \right]. \end{aligned}$$

□

9.2.2 Quadratic loss *

For quadratic loss ($L_0(a_1, a_2) := a_1^2$), the definition of $T^*(\mathbf{Y})$ in the above theorem is

$$T^*(\mathbf{Y}) = \arg \min_{a_1 \in \mathbb{R}} E \left[\left(d_1(\varepsilon) + d_2(\varepsilon)a_1 \right)^2 \middle| \mathbf{Y} \right].$$

We then have:

Lemma 9.2.1 Suppose that d is equivariant, and sufficient and complete. Then

$$T^*(\mathbf{X}) := d_1(\mathbf{X}) - d_2(\mathbf{X}) \frac{Ed_1(\varepsilon)d_2(\varepsilon)}{Ed_2^2(\varepsilon)}$$

is UMRE.

Proof. By Basu's lemma, d and \mathbf{Y} are independent. Hence

$$E\left[\left(d_1(\varepsilon) + d_2(\varepsilon)a_1\right)^2 \middle| \mathbf{Y}\right] = E\left(d_1(\varepsilon) + d_2(\varepsilon)a_1\right)^2.$$

Moreover

$$\arg \min_{a_1 \in \mathbb{R}} E\left(d_1(\varepsilon) + d_2(\varepsilon)a_1\right)^2 = -\frac{Ed_1(\varepsilon)d_2(\varepsilon)}{Ed_2^2(\varepsilon)}.$$

□

Example 9.2.1 UMRE of the mean of the normal distribution: σ^2 unknown

Let X_1, \dots, X_n be i.i.d. and $\mathcal{N}(\mu, \sigma^2)$ -distributed. Define

$$d_1(\mathbf{X}) := \bar{X}, \quad d_2(\mathbf{X}) := S,$$

where S^2 is the sample variance

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

It is easy to see that d is equivariant. We moreover know from Example 4.3.6 that d is sufficient, and an application of Lemma 5.4.1 shows that d is also complete. We furthermore have

$$Ed_1(\varepsilon) = E\bar{\epsilon} = 0,$$

and, from Example 9.1.2 (a consequence of Basu's lemma), we know that $d_1(\mathbf{X}) = \bar{X}$ and $d_2(\mathbf{X}) = S$ are independent. So

$$Ed_1(\varepsilon)d_2(\varepsilon) = Ed_1(\varepsilon)Ed_2(\varepsilon) = 0.$$

It follows that $T^*(\mathbf{X}) = \bar{X}$ is UMRE.

