

PIOTR ZWIERNIK

ADVANCED THEORY OF STATISTICS

Introduction

These lecture notes are intended to offer complement to the first part of STA3000. In designing this course we aimed at giving a modern treatment of mathematical statistics. Thus, we discuss important topics in theoretical statistics, like exponential families, statistical decision theory, or empirical processes and develop basic intuition behind what is multivariate statistics and what are its basic problems and techniques.

These goals heavily affected the exposition. Our approach is to use the multivariate notation whenever possible, to emphasize connections to convex analysis, and to present some results in the high-dimensional statistics. We will try to show that importance of convex analysis for statistical theory goes much beyond convex optimization used for the maximum likelihood estimation or its regularized versions.

The whole material is divided into twelve 3-hour lectures. The notes contain more detailed material than presented in the lecture. Preparing these lecture notes we benefited from several excellent textbooks or lecture notes:

1. Robert W. Keener, *Theoretical statistics*.
2. Lehmann, Romano, *Testing Statistical Hypotheses*.
3. Sundberg, *Statistical Modelling by Exponential Families*.
4. Wainwright, *High-dimensional statistics*.
5. Martin, *Lecture Notes on Advanced Statistical Theory*.
6. Rigolett, *High-dimensional Statistics*.
7. van der Vaart, *Asymptotic Statistics*

Special thanks go to Morris Greenberg, Ichiro Hashimoto, Vishakh Patel, Emily Somerset, Qiang Sun, Leonard Wang, and Zhenghang Xu for helping me to improve the notes.

Contents

	<i>I Topics in Statistical Inference</i>	9
1	<i>Exponential families (2 weeks)</i>	11
	1.1 <i>Definition and examples</i>	11
	1.2 <i>Basic results</i>	14
	1.3 <i>Convexity and the MLE</i>	17
	1.4 <i>Marginal and conditional distributions*</i>	22
	1.5 <i>Conditional inference for canonical parameter*</i>	25
	1.6 <i>Kullback-Leibler divergence</i>	26
	1.7 <i>Generalized Linear Models</i>	29
	1.8 <i>Diaconis-Ylvisaker conjugate priors</i>	31
	1.9 <i>Exercises</i>	32
2	<i>Statistical Decision Theory (1 week)</i>	35
	2.1 <i>Decision theoretic framework*</i>	35
	2.2 <i>Randomized decision rules*</i>	38
	2.3 <i>Bayesian and minimax rules*</i>	40
	2.4 <i>Admissibility and Rao-Blackwell</i>	47
	2.5 <i>Stein's paradox</i>	49
	2.6 <i>Minimizing risk under constraints</i>	54
	2.7 <i>Exercises</i>	62

3	<i>Hypothesis testing and multiple testing (2 weeks)</i>	67
3.1	<i>Neyman-Pearson lemma*</i>	68
3.2	<i>Some constructions of UMP tests*</i>	73
3.3	<i>Sequential testing*</i>	76
3.4	<i>Motivating multiple testing</i>	81
3.5	<i>Family-wise error rate</i>	81
3.6	<i>False discovery rate</i>	86
3.7	<i>Exercises</i>	90
	 <i>II Statistical Learning Theory: An Empirical Process Perspective</i>	 93
4	<i>Motivation and examples (1 week)</i>	95
4.1	<i>Uniform law of large numbers</i>	95
4.2	<i>The Uniform Central Limit Theory</i>	100
4.3	<i>Exercises</i>	103
5	<i>Concentration of measure (3 weeks)</i>	105
5.1	<i>Basic inequalities</i>	105
5.2	<i>Sub-gaussian and sub-exponential random variables</i>	109
5.3	<i>Martingale-based methods</i>	114
5.4	<i>Lipschitz functions of Gaussian variables</i>	118
5.5	<i>Exercises</i>	121
6	<i>More advanced techniques (1-2 weeks)</i>	123
6.1	<i>Maximal inequalities</i>	123
6.2	<i>Rademacher complexity and bounds on suprema</i>	126
6.3	<i>Polynomial discrimination and VC dimension</i>	128
6.4	<i>Chaining</i>	132
6.5	<i>Exercises</i>	135

7	<i>Applications in statistics (1-2 weeks)</i>	137
7.1	<i>Sub-Gaussian sequence model with sparsity</i>	137
7.2	<i>Fixed design linear regression</i>	140
7.3	<i>Constrained least squares estimator</i>	142
7.4	<i>LASSO regression</i>	143
7.5	<i>Random matrices</i>	147
7.6	<i>Exercises</i>	152
	<i>III Some topics from the first semester</i>	155
8	<i>Classical large sample theory</i>	157
8.1	<i>Preliminaries</i>	157
8.2	<i>Delta Method</i>	160
8.3	<i>M-estimators</i>	162
8.4	<i>Generalized likelihood ratio test</i>	165
8.5	<i>Limits of Bayesian procedures</i>	168
8.6	<i>Exercises</i>	170
	<i>IV Mathematical appendix</i>	171
A	<i>Real Analysis</i>	173
A.1	<i>Vector spaces</i>	173
A.2	<i>Continuity and semicontinuity</i>	175
A.3	<i>Differentiation</i>	177
B	<i>Convex Analysis</i>	181
B.1	<i>Convexity and hyperplane separation</i>	181
B.2	<i>Convex functions and optimization</i>	184
C	<i>Probability</i>	189
C.1	<i>Continuity of probability</i>	189
C.2	<i>Martingales</i>	190

Part I

**Topics in Statistical
Inference**

1

Exponential families (2 weeks)

Exponential families were discussed briefly in the first semester. The goal of this section is to provide a more detailed treatment of multivariate exponential families in connection with convexity, sufficiency, and hypothesis testing. This chapter is mostly based on two books on exponential families by Lawrence D. Brown¹ and by Rolf Sundberg², and on the lecture notes „Topics in Information Geometry“ of our colleague Ting-Kam Leonard Wong.

1.1 Definition and examples

Basic definition and univariate examples appeared in the first part of the lecture. We focus on developing uniform notation in the multivariate case. Consider a random vector $\mathbf{X} = (X_1, \dots, X_m)$ with values in the state space $\mathcal{X} \subseteq \mathbb{R}^m$ equipped with a σ -finite measure μ .³

Definition 1.1.1. *A parametric statistical model for \mathbf{X} is an exponential family with canonical parameter vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$ and canonical statistics $\mathbf{t}(\mathbf{x}) = (t_1(\mathbf{x}), \dots, t_d(\mathbf{x}))$, if it admits a density f with respect to μ and f has the form*

$$f(\mathbf{x}; \boldsymbol{\theta}) = h(\mathbf{x}) \exp \left\{ \langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle - A(\boldsymbol{\theta}) \right\}. \quad (1.1)$$

Formally, we define \mathcal{X} as the smallest closed set satisfying $\mathbb{P}_\theta(X \in \mathcal{X}) = 1$. This definition does not depend on the choice of $\boldsymbol{\theta}$ because all \mathbb{P}_θ have the same support.

Remark 1.1.2. *For notational purposes it is often easier to subsume $h(\mathbf{x})$ into the underlying measure and use the formulation*

$$f(\mathbf{x}; \boldsymbol{\theta}) = \exp \left\{ \langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle - A(\boldsymbol{\theta}) \right\}. \quad (1.2)$$

Another useful reformulation is when $h(\mathbf{x})$ in (1.1) is itself a density. In this case $A(\mathbf{0}) = 0$.

¹ Lawrence D. Brown. *Fundamentals of statistical exponential families with applications in statistical decision theory*, volume 9 of *Institute of Mathematical Statistics Lecture Notes—Monograph Series*. Institute of Mathematical Statistics, Hayward, CA, 1986

² Rolf Sundberg. *Statistical modelling by exponential families*, volume 12 of *Institute of Mathematical Statistics Textbooks*. Cambridge University Press, Cambridge, 2019

³ Typically the measure μ is either the counting measure or the Lebesgue measure.

Example 1.1.3 (Bernoulli variable and logistic regression). If $X \sim \text{Bern}(p)$ with $p \in (0, 1)$ we can write its distribution in the exponential form. For $x \in \{0, 1\}$ we have

$$f(x; p) = p^x(1-p)^{1-x} = \exp \left\{ x \log \left(\frac{p}{1-p} \right) + \log(1-p) \right\}.$$

We have $t(x) = x$ and the canonical parameter is the logit of p :

$$\text{logit}(p) := \log \frac{p}{1-p} = \theta.$$

We also have $h(x) = 1$ and $A(\theta) = \log(1 + e^\theta)$. In logistic regression, we model $\text{logit}(p)$ as a linear function of regressors.

Example 1.1.4 (Univariate Gaussian). If $X \sim N(\mu, \sigma^2)$ then

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} + \frac{1}{2} \left(\log\left(\frac{1}{\sigma^2}\right) - \frac{\mu^2}{\sigma^2} \right)},$$

which can be written as a two dimensional exponential family with

$$\mathbf{t}(x) = \left(x, -\frac{x^2}{2} \right), \quad \boldsymbol{\theta} = \left(\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2} \right).$$

Then

$$h(\mathbf{x}) = \frac{1}{\sqrt{2\pi}}, \quad A(\boldsymbol{\theta}) = -\frac{1}{2} \left(\log(\theta_2) - \frac{\theta_1^2}{\theta_2} \right).$$

If $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ are i.i.d. then the joint distribution of this sample $\mathbf{X}_{1:n} = (X_1, \dots, X_n)$ is also Gaussian with the same canonical parameters and the density

$$f(\mathbf{x}_{1:n}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum_i x_i^2 + \frac{\mu}{\sigma^2} \sum_i x_i + \frac{n}{2} \left(\log\left(\frac{1}{\sigma^2}\right) - \frac{\mu^2}{\sigma^2} \right)}$$

and so it forms an exponential family with the same canonical parameter and with the sufficient statistics $\left(\sum_i x_i, -\frac{1}{2} \sum_i x_i^2 \right)$. We know, of course, that this distribution is the n -variate Gaussian with parameters $\mu \mathbf{1}$ and $\sigma^2 I_n$.

Consider the function

$$Z(\boldsymbol{\theta}) := \int_{\mathcal{X}} h(\mathbf{x}) \exp \left\{ \langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle \right\} \mu(d\mathbf{x}), \quad (1.3)$$

where we put $Z(\boldsymbol{\theta}) = +\infty$ if this integral is infinite. Because f is a density function, it follows that

$$A(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta}).$$

The function $A(\boldsymbol{\theta})$ plays a special role in this theory and it has many names: the log-partition function, Laplace transform, or the cumulant function. We define the space of canonical parameters as

$$\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^d : Z(\boldsymbol{\theta}) < +\infty\}. \quad (1.4)$$

It is implicit in (1.1) that $\boldsymbol{\theta} \in \Theta$. An exponential family model is then specified by possibly constraining to $\Theta_0 \subseteq \Theta$. The dimension d of $\boldsymbol{\theta}$ is called the order of this exponential family.

Example 1.1.5. Consider the univariate Gaussian case in Example 1.1.4. Since $\mu \in \mathbb{R}$, $\sigma^2 > 0$, the space of canonical parameters is $\Theta = \{(\theta_1, \theta_2) : \theta_1 \in \mathbb{R}, \theta_2 > 0\}$. Taking $\mu = 0$ corresponds to fixing a linear subspace $\Theta_0 = \{\boldsymbol{\theta} \in \Theta : \theta_1 = 0\}$.

In our basic set-up both the parameters and the sufficient statistics lie in \mathbb{R}^d but exactly the same definition can be provided for a general d -dimensional vector space with a given inner product. The most relevant example is when the underlying vector space is the space \mathbb{S}^m of all symmetric $m \times m$ matrices. Here the standard inner product is given by $\langle A, B \rangle = \text{tr}(AB)$. Denote by \mathbb{S}_+^m the set of positive definite matrices in \mathbb{S}_+^m .

Show that $\text{tr}(AB) = \sum_{ij} A_{ij}B_{ij}$

Example 1.1.6 (Centered multivariate Gaussian distribution). Consider the m -variate Gaussian distribution with the zero mean vector and covariance matrix $\Sigma \in \mathbb{S}_+^m$. Let $K = \Sigma^{-1}$. The density with respect to the Lebesgue measure is

$$f(\mathbf{x}; K) = \frac{1}{(2\pi)^{m/2}} \sqrt{\det(K)} \exp\{-\frac{1}{2}\mathbf{x}^\top K \mathbf{x}\}.$$

This is an exponential family with $h(\mathbf{x}) = \frac{1}{(2\pi)^{m/2}}$, $A(K) = -\frac{1}{2} \log \det(K)$. Denoting $\mathbf{t}(\mathbf{x}) = -\frac{1}{2}\mathbf{x}\mathbf{x}^\top \in \mathbb{S}^m$ we get

$$-\frac{1}{2}\mathbf{x}^\top K \mathbf{x} = \text{tr}(K\mathbf{t}(\mathbf{x})) = \langle K, \mathbf{t}(\mathbf{x}) \rangle.$$

← Exercise 1.9.1

← Exercise 1.9.2

Given a sample $\mathbf{x}_{1:n} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ from a distribution with density $f(\mathbf{x}; \boldsymbol{\theta})$, the log-likelihood function is

$$\ell_n(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}_i; \boldsymbol{\theta}).$$

Note that we normalize the log-likelihood by n to get the interpretation as the expectation of $\log f(\mathbf{x}; \boldsymbol{\theta})$ under the sample distribution.

Proposition 1.1.7 (The log-likelihood function). If $\mathbf{x}_{1:n} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ is a sample from the exponential family (1.1) then, denoting

← Exercise 1.9.3

$$\bar{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{t}(\mathbf{x}_i),$$

the log-likelihood takes the form

$$\ell_n(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \bar{\boldsymbol{\mu}}_n \rangle - A(\boldsymbol{\theta}) + (\text{constant}). \quad (1.5)$$

In the statistical/ML practice we typically go the other way around. We first define a suitable sufficient statistics that should contain all the relevant information of the data. This choice defines then an exponential family.

Example 1.1.8 (Exponential random graph model). *Each graph can be associated to its adjacency matrix $A \in \{0, 1\}^{n \times n}$. This is a symmetric matrix with zeros on the diagonal and so each graph is an element $\mathbf{x} \in \mathcal{X} = \{0, 1\}^{\binom{n}{2}}$ with entries x_{ij} for $1 \leq i < j \leq n$. For the simplest example, consider $\mathbf{t}(\mathbf{x}) = \sum_{i < j} x_{ij} \in \mathbb{R}$, which is simply the number of edges of the underlying graph. The corresponding exponential family has one parameter θ and is of the form*

$$f(\mathbf{x}; \boldsymbol{\theta}) = \exp \left\{ \theta \sum_{i < j} x_{ij} - A(\theta) \right\} \propto \prod_{i < j} e^{\theta x_{ij}}.$$

In other words, each edge x_{ij} is an independent Bernoulli variable with the success probability $p = e^\theta / (1 + e^\theta)$ and computing the normalizing constant is easy. This is the famous Erdős-Renyi model. Other statistics of the graph will give different exponential models.

Another model for binary variables that uses the graph structure is the Ising model.

Example 1.1.9 (The Ising model). *Let G be a graph over m nodes representing m binary random variables $X_i \in \{-1, 1\}$ for $i = 1, \dots, m$. Consider the model⁴*

$$f(\mathbf{x}; \boldsymbol{\theta}) \propto \exp \left\{ \sum_{ij \in G} \theta_{ij} x_i x_j \right\}.$$

Here computing the normalizing constant is generally hard. This and similar examples motivated developing methods that do not rely on computing this normalizing constant (e.g. variational inference).

There are many models with a given sufficient statistics. What is then special about the exponential families? To answer this question we need to develop more theory and we will see that this is connected to maximizing the entropy; c.f. Theorem 1.6.5.

1.2 Basic results

We assume throughout that there is no hyperplane $H = \{\mathbf{x} : \langle \boldsymbol{\alpha}, \mathbf{t}(\mathbf{x}) \rangle = c\}$ such that $\mathbb{P}_\theta(H) = 1$. In other words, no entry of the vector \mathbf{t} can be written as an affine combination of the remaining entries. Similarly, we assume that there is no hyperplane

⁴ The statistical interpretation of this modelling construction in terms of conditional independence comes from the Hammersley-Clifford theorem.

← Exercise 1.9.4

← Exercise 1.9.5

For example $\langle (1, 2), (x_1, x_2) \rangle = x_1 + 2x_2 = 1$ is a simple hyperplane in \mathbb{R}^2 .

$\{\boldsymbol{\theta} : \langle \boldsymbol{\theta}, \boldsymbol{\beta} \rangle = c\}$ containing Θ . An exponential family (1.1) with these two properties is called **minimal**. A canonical example of a family that is not minimal appears for discrete data.

Example 1.2.1 (Bernoulli variable). *The Bernoulli distribution in Example 1.1.3 could be alternatively written in form of a two-dimensional exponential family*

$$p(x) = \exp\{\log(1-p)\mathbb{1}(x=0) + \log(p)\mathbb{1}(x=1)\} \quad \text{for } x \in \{0,1\}$$

with $\mathbf{t}(x) = (\mathbb{1}(x=0), \mathbb{1}(x=1))$, $\boldsymbol{\theta} = (\log(1-p), \log(p))$, $A(\boldsymbol{\theta}) = 0$.

Clearly, this representation is not minimal as $\mathbb{1}(x=0) + \mathbb{1}(x=1) = 1$.

An easy fix is to define the new (minimal) sufficient statistics $x = \mathbb{1}(x=1)$ and rewrite the above using the fact that $\mathbb{1}(x=0) = 1-x$. This gives the representation in Example 1.1.3.

A slightly more complicated version of this example appears in the vector case. But the idea is similar and the difficulty is purely notational.

Example 1.2.2 (Binary vectors). *Consider a binary vector $X = (X_1, \dots, X_m)$ with the probability distribution $p(\mathbf{x})$ for $\mathbf{x} \in \mathcal{X} = \{0,1\}^m$. As in the Bernoulli case, we can write*

$$p(\mathbf{x}) = \exp\{\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle\} \quad \text{for } \mathbf{x} \in \mathcal{X},$$

where $\boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \in \mathbb{R}^{\mathcal{X}}$. By this, we mean that $\boldsymbol{\theta}$ and $\mathbf{t}(\mathbf{x})$ are themselves functions on \mathcal{X} : $\boldsymbol{\theta}(\mathbf{y})$ and $\mathbf{t}(\mathbf{x}, \mathbf{y})$ for $\mathbf{y} \in \mathcal{X}$, such that $\boldsymbol{\theta}(\mathbf{y}) = \log p(\mathbf{y})$, and

$$\mathbf{t}(\mathbf{x}, \mathbf{y}) = \mathbb{1}(\mathbf{x} = \mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{y}, \\ 0 & \text{otherwise.} \end{cases}$$

The inner product in $\mathbb{R}^{\mathcal{X}}$ simply means

$$\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle = \sum_{\mathbf{y} \in \mathcal{X}} \boldsymbol{\theta}(\mathbf{y}) \mathbf{t}(\mathbf{x}, \mathbf{y})$$

However $\sum_{\mathbf{x} \in \mathcal{X}} \mathbf{t}(\mathbf{x}) = \mathbf{1} \in \mathbb{R}^{\mathcal{X}}$ so this representation is not minimal. We reduce the dimension by rewriting

$$\mathbf{t}(\mathbf{0}) = \mathbf{1} - \sum_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\}} \mathbf{t}(\mathbf{x}).$$

Show that this defines a minimal exponential family with $\Theta = \mathbb{R}^{\mathcal{X} \setminus \{\mathbf{0}\}}$ and canonical parameters $\boldsymbol{\theta}(\mathbf{y}) = \log p(\mathbf{y}) - \log p(\mathbf{0})$ and

$$A(\boldsymbol{\theta}) = \log \left(1 + \sum_{\mathbf{x} \neq \mathbf{0}} e^{\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle_0} \right),$$

where the inner product $\langle \cdot, \cdot \rangle_0$ is defined in the space $\mathbb{R}^{\mathcal{X} \setminus \{\mathbf{0}\}}$, that is,

$$\langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle_0 = \sum_{\mathbf{y} \neq \mathbf{0}} \boldsymbol{\theta}(\mathbf{y}) \mathbf{t}(\mathbf{x}, \mathbf{y}).$$

Making sure that there is no hyperplane containing $t(\mathcal{X})$ can be easily done by changing the sufficient statistics as in the examples above⁵. By redefining Θ if necessary, we can also always without loss assume that the canonical parameters are not contained in an affine subspace.

Example 1.2.3 (Gaussian graphical models). Consider the centered m -variate Gaussian model in Example 1.1.6 with the canonical parameter $K = \Sigma^{-1} \in \Theta = \mathbb{S}_+^m$. Fix a graph G over m nodes and the subset

$$\Theta_G = \{K \in \Theta : K_{ij} = 0 \text{ if } ij \notin G\}.$$

This linear constrain defines an exponential family with sufficient statistic $(-\frac{1}{2}x_i x_j)_{ij \in G}$. The corresponding model, is typically called the Gaussian graphical model, which forms a popular dimension reduction technique. Its relevance in practice comes from the fact that in the Gaussian distribution $K_{ij} = 0$ if and only X_i is independent of X_j given all the remaining variables in the system.⁶

Remark 1.2.4. The rest of this section is not essential if you skip Section 1.4 and Section 1.5.

Definition 1.2.5. A minimal exponential family is called **full** if its parameter space is maximal, that is, Θ_0 equals the canonical space Θ .

Some relevant examples of a non-full exponential family are when the parameter space Θ_0 is a convex subset of Θ (convex exponential families) or when it forms a lower dimensional manifold (curved exponential families). Note however that if this manifold is a linear subspace we again get a full exponential family after a reparametrization.

Recall from the first semester that a statistic $\mathbf{t}(\mathbf{x})$ is **sufficient** for θ if the conditional distribution of \mathbf{x} given $\mathbf{t}(\mathbf{x})$ does not depend on θ . Sufficiency of $\mathbf{t}(\mathbf{x})$ in exponential families can be argued by Proposition 1.2.7 below or directly by Fisher-Neyman Factorization Theorem, which states that $\mathbf{t}(\mathbf{x})$ is sufficient for θ if and only if $p(\mathbf{x}; \theta) = h(\mathbf{x})g_\theta(\mathbf{t}(\mathbf{x}))$ for some h, g_θ . Typically the canonical statistics is also minimal sufficient, that is, for any other sufficient statistic $\mathbf{t}'(\mathbf{x})$,

$$\mathbf{t}'(\mathbf{x}) = \mathbf{t}'(\mathbf{y}) \quad \implies \quad \mathbf{t}(\mathbf{x}) = \mathbf{t}(\mathbf{y}).$$

Proposition 1.2.6 (Minimal sufficiency of \mathbf{t}). In a full exponential family the statistic $\mathbf{t}(\mathbf{x})$ is minimally sufficient for θ .

Proof. Consider any other sufficient statistics \mathbf{t}' . If $\mathbf{t}'(\mathbf{x}) = \mathbf{t}'(\mathbf{y})$ then the factorization theorem shows that $\frac{f(\mathbf{x}; \theta)}{f(\mathbf{y}; \theta)}$ is independent of θ . On the other hand

$$\frac{f(\mathbf{x}; \theta)}{f(\mathbf{y}; \theta)} = \frac{h(\mathbf{x})}{h(\mathbf{y})} \exp\{\langle \theta, \mathbf{t}(\mathbf{x}) - \mathbf{t}(\mathbf{y}) \rangle\},$$

⁵ If one coordinate of $t(\mathcal{X})$ can be written as an affine combination of the others, we simply replace this coordinate with this combination obtaining a sufficient statistics with one dimension lower. This process can be repeated.

← Exercise 1.9.6

← Exercise 1.9.5

⁶ Caroline Uhler. Gaussian graphical models. In *Handbook of Graphical Models*, pages 217–238. CRC Press, 2018

← Exercise 1.9.6

which can be constant in θ if and only if $\langle \theta, \mathbf{t}(\mathbf{x}) - \mathbf{t}(\mathbf{y}) \rangle$ is constant. Because the family is full (in particular θ is not contained in an affine subspace), this happens if and only if $\mathbf{t}(\mathbf{x}) = \mathbf{t}(\mathbf{y})$. \square

Note that in the proof we only used that θ is not contained in an affine space, so this result generalizes to curved exponential families.

Proposition 1.2.7 (Distribution of the sufficient statistic). *Suppose \mathbf{X} has distribution in the exponential family (1.1). Then, under certain regularity conditions, the distribution of $\mathbf{t} = \mathbf{t}(\mathbf{X})$ is*

$$f(\mathbf{t}; \theta) = g(\mathbf{t}) \exp\{\langle \theta, \mathbf{t} \rangle - A(\theta)\}, \quad (1.6)$$

where the structure function $g(\mathbf{t})$ in the discrete case is

$$g(\mathbf{t}) = \sum_{\mathbf{t}(\mathbf{x})=\mathbf{t}} h(\mathbf{x}),$$

and in the continuous case

$$g(\mathbf{t}) = \int_{\mathbf{t}(\mathbf{x})=\mathbf{t}} h(\mathbf{x}) d\mathbf{x}.$$

Example 1.2.8. Consider a centered Gaussian distribution with $\Sigma \in \mathbb{S}_+^m$. Given a random sample $\mathbf{x}_1, \dots, \mathbf{x}_n$ from this distribution, the statistics $n\bar{\boldsymbol{\mu}}_n := \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ has the Wishart distribution.

In the discrete case the proof is elementary.

The last expression is written rather informally. The integral is computed with respect to the measure on the set $\{\mathbf{x} : \mathbf{t}(\mathbf{x}) = \mathbf{t}\}$ induced from the Lebesgue measure. Here the proof is non-trivial and we skip it.

1.3 Convexity and the MLE

Recall a basic version of the Hölder's inequality.

Proposition 1.3.1 (Hölder's inequality). *If f, g are two functions on a measurable space (\mathcal{X}, μ) then for every $p, q \in [1, \infty]$ such that $1/p + 1/q = 1$ we have*

$$\int_{\mathcal{X}} |f(\mathbf{x})g(\mathbf{x})| \mu(d\mathbf{x}) \leq \left(\int_{\mathcal{X}} |f(\mathbf{x})|^p \mu(d\mathbf{x}) \right)^{1/p} \left(\int_{\mathcal{X}} |g(\mathbf{x})|^q \mu(d\mathbf{x}) \right)^{1/q}. \quad (1.7)$$

Moreover, if $p, q > 1$ then (1.7) holds as equality if and only if $|f|^p$ and $|g|^q$ are linearly dependent in $L^1(\mathcal{X})$ meaning that there exist real numbers $\alpha, \beta \geq 0$ such that $\alpha|f|^p = \beta|g|^q$ μ -almost everywhere.

$L^1(\mathcal{X})$ is the vector space of all functions f on \mathcal{X} with the property that $\|f\|_1 := \int_{\mathcal{X}} |f(x)| \mu(d\mathbf{x})$ is finite. Formally, we identify two functions that are equal almost surely.

Here is an important fundamental fact about exponential families.

Theorem 1.3.2. *For every exponential family (1.1) with Θ defined in (1.4) we have:*

- (i) Θ is a convex set and the function $A(\theta)$ is convex on Θ .
- (ii) $\mathbb{P}_{\theta_1} = \mathbb{P}_{\theta_2}$ if and only if

$$A((1-\lambda)\theta_1 + \lambda\theta_2) = (1-\lambda)A(\theta_1) + \lambda A(\theta_2) \quad \text{for all } \lambda \in (0, 1).$$

(iii) If the exponential family is minimal then A is strictly convex on Θ and $\mathbb{P}_{\theta_1} \neq \mathbb{P}_{\theta_2}$ if $\theta_1 \neq \theta_2 \in \Theta$.

(iv) A is lower semi-continuous on \mathbb{R}^d and is continuous in the interior of Θ .

For the definitions and basic results see Section A.2.2.

Proof. (i) Let $\theta_1, \theta_2 \in \Theta$, $\lambda \in (0, 1)$ and denote $\theta_\lambda = (1 - \lambda)\theta_1 + \lambda\theta_2$. By the Hölder's inequality with $p = 1/(1 - \lambda)$ and $q = 1/\lambda$:

$$\begin{aligned} Z(\theta_\lambda) &= \int h(\mathbf{x}) e^{\langle \theta_\lambda, \mathbf{t}(\mathbf{x}) \rangle} \mu(d\mathbf{x}) \\ &= \int \left(h(\mathbf{x}) e^{\langle \theta_1, \mathbf{t}(\mathbf{x}) \rangle} \right)^{1-\lambda} \left(h(\mathbf{x}) e^{\langle \theta_2, \mathbf{t}(\mathbf{x}) \rangle} \right)^\lambda \mu(d\mathbf{x}) \\ &\leq \left(\int h(\mathbf{x}) e^{\langle \theta_1, \mathbf{t}(\mathbf{x}) \rangle} \mu(d\mathbf{x}) \right)^{1-\lambda} \left(\int h(\mathbf{x}) e^{\langle \theta_2, \mathbf{t}(\mathbf{x}) \rangle} \mu(d\mathbf{x}) \right)^\lambda \\ &= Z(\theta_1)^{1-\lambda} Z(\theta_2)^\lambda. \end{aligned}$$

Taking the logs we get convexity of A . Now convexity of Θ follows easily.

(ii) The Hölder's inequality above is strict unless

$$\langle \theta_1 - \theta_2, \mathbf{t}(\mathbf{x}) \rangle \equiv \text{const} \quad (\mu \text{ a.s.}). \quad (1.8)$$

This last assertion is equivalent to $\mathbb{P}_{\theta_1} = \mathbb{P}_{\theta_2}$.

(iii) If (1.8) holds for some $\theta_1 \neq \theta_2$ then the exponential family is not minimal.

(iv) By Fatou's lemma $Z(\theta)$ is lower semicontinuous and so $A(\theta)$ is lower semicontinuous. Any convex function defined and finite on a convex set $\Theta \subset \mathbb{R}^d$ must be continuous on the interior of Θ . \square

← Exercise B.2.10

The following definition will be important in the rest of this chapter.

Definition 1.3.3. A minimal exponential family is **regular** if its canonical parameter space Θ is open.

In the finite discrete case we are always regular.

← Exercise 1.9.7

Most of the exponential families you will ever encounter are regular. But Θ does not need to be always open. An instance of a non-regular exponential family is given by the inverse Gaussian⁷.

⁷ See Section 3.2.1 in Rolf Sundberg's book

Proposition 1.3.4. In a regular exponential family, $A(\theta)$ is smooth and

$$\nabla A(\theta) = \mathbb{E}_\theta(\mathbf{t}(\mathbf{X})) =: \mu(\theta) \quad (1.9)$$

$$\nabla^2 A(\theta) = \text{var}_\theta(\mathbf{t}(\mathbf{X})) =: V(\theta). \quad (1.10)$$

Taking higher derivatives, we obtain higher cumulants of $\mathbf{t}(\mathbf{X})$.

Proof. Note that

$$K(\mathbf{s}) := A(\boldsymbol{\theta}_0 + \mathbf{s}) - A(\boldsymbol{\theta}_0) = \log \mathbb{E}_{\boldsymbol{\theta}_0} e^{\langle \mathbf{s}, \mathbf{t}(\mathbf{x}) \rangle},$$

which shows that $K(\mathbf{s})$ is the cumulant generating function of $\mathbf{t}(X)$ with respect to the distribution $\mathbb{P}_{\boldsymbol{\theta}_0}$. Since $\boldsymbol{\theta}_0$ is an interior point of Θ , $K(\mathbf{s})$ is well-defined in a neighbourhood of zero and so all cumulants exist (see, for example, p. 267 in ⁸). Using the chain rule, we get that, for any $r \geq 1$, the r -th order derivatives satisfy $\nabla_{\mathbf{s}}^r K(\mathbf{s}) = \nabla_{\boldsymbol{\theta}}^r A(\boldsymbol{\theta}_0 + \mathbf{s})$, where the notation ∇^r is hopefully intuitive and it denotes a $d \times \dots \times d$ array with the (i_1, \dots, i_r) -th entry:

$$(\nabla_{\mathbf{s}}^r K(\mathbf{0}))_{i_1 \dots i_r} = \frac{\partial^r}{\partial s_{i_1} \dots \partial s_{i_r}} K(\mathbf{0}), \quad (\nabla_{\boldsymbol{\theta}}^r A(\boldsymbol{\theta}_0))_{i_1 \dots i_r} = \frac{\partial^r}{\partial \theta_{i_1} \dots \partial \theta_{i_r}} A(\boldsymbol{\theta}_0).$$

By taking $\mathbf{s} = \mathbf{0}$, we obtain

$$\nabla_{\mathbf{s}}^r K(\mathbf{0}) = \nabla_{\boldsymbol{\theta}}^r A(\boldsymbol{\theta}_0), \quad (1.11)$$

The expression on the left in (1.11) gives precisely the r -th order cumulants of $\mathbf{t}(X)$. If $r = 1$ this is $\mathbb{E}_{\boldsymbol{\theta}_0}(\mathbf{t}(X))$ and if $r = 2$ then this is $\text{var}_{\boldsymbol{\theta}_0}(\mathbf{t}(X))$. \square

If $A(\boldsymbol{\theta})$ is smooth then $\log f(\mathbf{x}; \boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{t}(\mathbf{x}) \rangle - A(\boldsymbol{\theta})$ is smooth too. In this case the **Fisher information matrix** satisfies

$$I(\boldsymbol{\theta}) := -\mathbb{E}_{\boldsymbol{\theta}}(\nabla_{\boldsymbol{\theta}}^2 \log f(\mathbf{x}; \boldsymbol{\theta})) = \nabla_{\boldsymbol{\theta}}^2 A(\boldsymbol{\theta}).$$

Note that the observed information

$$J(\boldsymbol{\theta}) := -\nabla_{\boldsymbol{\theta}}^2 \ell_n(\boldsymbol{\theta}) = I(\boldsymbol{\theta})$$

and so, in particular, it does not depend on the data. We obtain the following result.

Proposition 1.3.5. *In a regular exponential family, the log-likelihood function, given in (1.5), is a smooth and strictly concave function of the canonical parameter $\boldsymbol{\theta}$. The score function $U(\boldsymbol{\theta}) = \nabla \ell_n(\boldsymbol{\theta})$ satisfies*

$$U(\boldsymbol{\theta}) = \bar{\boldsymbol{\mu}}_n - \boldsymbol{\mu}(\boldsymbol{\theta})$$

and the observed information $J(\boldsymbol{\theta})$ equals the expected (Fisher) information $I(\boldsymbol{\theta})$, and they are both given by the variance of \mathbf{t} ,

$$I(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) = V(\boldsymbol{\theta}).$$

Recall that \mathcal{X} is formally defined as the support of $\mathbb{P}_{\boldsymbol{\theta}}$. Define \mathcal{K} to be the convex hull of the image $\mathbf{t}(\mathcal{X})$:

$$\mathcal{K} := \text{conv}(\mathbf{t}(\mathcal{X})). \quad (1.12)$$

⁸ Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013

For example, for the multivariate Gaussian distribution in Example 1.1.6, $\mathbf{x} \in \mathbb{R}^m$ and so $\mathbf{t}(\mathbf{x}) = -\frac{1}{2}\mathbf{x}\mathbf{x}^\top$ is a rank-one negative semi-definite matrix. In this case, the convex hull of all such matrices, namely \mathcal{K} , is the cone of all negative semidefinite matrices.

The next results offers an alternative parametrization for the exponential model.

Proposition 1.3.6. *In a minimal regular exponential family:*

- (i) *The mapping $\mu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ given by $\boldsymbol{\theta} \mapsto \nabla A(\boldsymbol{\theta})$ is one-to-one on Θ .*
- (ii) *The log-likelihood function with data $\bar{\boldsymbol{\mu}} \in \mathcal{K}$ has a maximum in Θ if and only if $\bar{\boldsymbol{\mu}} \in \mu(\Theta)$ and then the maximum $\bar{\boldsymbol{\theta}}$ is given uniquely as $\bar{\boldsymbol{\theta}} = \mu^{-1}(\bar{\boldsymbol{\mu}})$.*
- (iii) *$\mu(\Theta) = \text{int}(\mathcal{K})$ and so in particular $\mu(\Theta)$ is open.*

Proof. (i) By Theorem 1.3.2, minimality implies that A is strictly convex. Thus, for every $\mathbf{m} \in \mathbb{R}^d$, the function $\langle \boldsymbol{\theta}, \mathbf{m} \rangle - A(\boldsymbol{\theta})$ is strictly concave. In particular, it has at most one stationary point in Θ , that is, at most one $\bar{\boldsymbol{\theta}} \in \Theta$ such $\mathbf{m} = \nabla A(\bar{\boldsymbol{\theta}})$.

(ii) This is just (i) rephrased.

(iii) To get the first inclusion $\mu(\Theta) \subseteq \text{int}(\mathcal{K})$, we first show that $\mu(\Theta) \subseteq \bar{\mathcal{K}}$, where the latter denotes the closure of \mathcal{K} . For every $\boldsymbol{\theta} \in \Theta$ and for every $c \in \mathbb{R}$, if there exists $\mathbf{u} \in \mathbb{R}^d$ such that $\langle \mathbf{u}, \mathbf{t}(\mathbf{X}) \rangle \leq c$ almost surely ($\mathbf{t}(\mathcal{X})$ is contained in the given half-space) then $\langle \mathbf{u}, \mathbb{E}_{\boldsymbol{\theta}}(\mathbf{t}(\mathbf{X})) \rangle \leq c$ for every $\boldsymbol{\theta} \in \Theta$. This implies that if a half-space H contains \mathcal{K} then it also contains $\mu(\Theta)$. The intersection of all such halfplanes is equal to $\bar{\mathcal{K}}$; this is a standard application of the Hyperplane Separation Theorem B.1.3. This shows that $\mu(\Theta) \subseteq \bar{\mathcal{K}}$. To prove that $\mu(\Theta) \subseteq \text{int}(\mathcal{K})$, we argue by contradiction. Suppose that $\mu(\boldsymbol{\theta}_0)$ for some $\boldsymbol{\theta}_0 \in \Theta$ lies in the boundary of \mathcal{K} , $\partial\mathcal{K} = \bar{\mathcal{K}} \setminus \text{int}(\mathcal{K})$. Then, again by the Hyperplane Separation Theorem, there exists a closed halfplane $H = \{\boldsymbol{\mu} : \langle \mathbf{u}, \boldsymbol{\mu} \rangle \leq c\}$ such that $\text{int}(\mathcal{K}) \subseteq H$ but $\mu(\boldsymbol{\theta}_0) \in \partial H$. In particular,

$$Z := \langle \mathbf{u}, \mathbf{t}(\mathbf{X}) - \mu(\boldsymbol{\theta}_0) \rangle \leq 0 \quad \text{almost surely.}$$

Note however that, since $\mathbb{E}_{\boldsymbol{\theta}} \mathbf{t}(\mathbf{X}) = \mu(\boldsymbol{\theta})$, $\mathbb{E}_{\boldsymbol{\theta}_0} Z = 0$, which implies that $Z = 0$ $\mathbb{P}_{\boldsymbol{\theta}_0}$ -almost surely. This however contradicts minimality of the exponential family. We conclude that $\mu(\Theta) \subseteq \text{int}(\mathcal{K})$.

It remains to show that the opposite inclusion $\text{int}(\mathcal{K}) \subseteq \mu(\Theta)$ also holds. We again argue by contradiction. Let $\mathbf{t}_0 \in \text{int}(\mathcal{K}) \setminus \mu(\Theta)$ then the equation $\mathbf{t}_0 = \nabla A(\boldsymbol{\theta})$ has no solution, or equivalently, the log-likelihood function $\ell(\boldsymbol{\theta}) = \langle \boldsymbol{\theta}, \mathbf{t}_0 \rangle - A(\boldsymbol{\theta})$ is not bounded above on Θ . To get a contradiction, fix any $\boldsymbol{\theta}_0 \in \Theta$ and consider half-lines

$$L_{\mathbf{u}} := \{\boldsymbol{\theta}_\lambda = \boldsymbol{\theta}_0 + \lambda \mathbf{u} : \lambda \geq 0\}.$$

← Exercise B.1.4

Now L_u can be either not entirely contained in Θ (case 1) or it can be contained in Θ (case 2). We will show that over each L_u the value of $-\ell(\theta_\lambda)$ goes to infinity as $\lambda \rightarrow \infty$ (in case 1) or as we approach the boundary of Θ (in case 2). In consequence, maximizing ℓ over Θ can be reduced to a compact subset $\{\theta : \ell(\theta) \geq \ell(\theta_0)\}$ and thus the optimum must exist. But then it must be of the form $\mathbf{t}_0 = \nabla A(\hat{\theta})$ leading to a contradiction.

Case 1: Since the half-line L_u is not contained in Θ , for some λ_0 the point θ_{λ_0} lies on the boundary of Θ . Since Θ is open, $A(\theta_{\lambda_0}) = \infty$ and hence $\ell(\theta_{\lambda_0}) = -\infty$. Thus $-\ell(\theta_\lambda) \rightarrow \infty$ as $\lambda \rightarrow \lambda_0$.

Case 2: In this case we can take $\lambda \rightarrow \infty$ without leaving Θ . We will still show that $-\ell(\theta_\lambda) \rightarrow \infty$. Note that

$$e^{-\ell(\theta_\lambda)} = Z(\theta_\lambda)e^{-\langle \theta_\lambda, \mathbf{t}_0 \rangle} = \int_{\mathcal{X}} h(\mathbf{x})e^{\lambda \langle \mathbf{u}, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle} e^{\langle \theta_0, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle} \mu(d\mathbf{x}).$$

Denote the integrand by $I(\lambda)$. Define

$$\begin{aligned} A_+ &= \{\mathbf{x} \in \mathcal{X} : \langle \mathbf{u}, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle > 0\}, \\ A_0 &= \{\mathbf{x} \in \mathcal{X} : \langle \mathbf{u}, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle = 0\}, \\ A_- &= \{\mathbf{x} \in \mathcal{X} : \langle \mathbf{u}, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle < 0\}, \end{aligned}$$

and note that A_+, A_0, A_- form a partition of \mathcal{X} . The integrand $I(\lambda)$ is always non-negative and it is increasing in λ on A_+ , constant on A_0 , and decreasing on A_- . The *monotone convergence theorem* assures that we can pass with the limit of λ inside the integral, namely

$$\lim_{\lambda \rightarrow \infty} \int_{A_-} I(\lambda) \mu(d\mathbf{x}) = 0,$$

$$\lim_{\lambda \rightarrow \infty} \int_{A_0} I(\lambda) \mu(d\mathbf{x}) = \int_{A_0} e^{\langle \theta_0, \mathbf{t}(\mathbf{x}) - \mathbf{t}_0 \rangle} \mu(d\mathbf{x}) < \infty,$$

and

$$\lim_{\lambda \rightarrow \infty} \int_{A_+} I(\lambda) \mu(d\mathbf{x}) = \infty.$$

Unless A_+ has measure zero, we conclude that $\lim_{\lambda \rightarrow \infty} e^{-\ell(\theta_\lambda)} = \infty$ or equivalently $\lim_{\lambda \rightarrow \infty} -\ell(\theta_\lambda) = \infty$.

To conclude the proof, it remains to show that A_+ has positive measure. If the measure is zero, $\langle \mathbf{u}, \mathbf{t}(\mathbf{x}) \rangle \leq \langle \mathbf{u}, \mathbf{t}_0 \rangle$ for all $\mathbf{x} \in \mathcal{X}$. We will again show that there must be equality, which contradicts minimality. This is where we use the fact that \mathbf{t}_0 is an interior point of \mathcal{K} . Let $\mathbf{x}_1 \in \mathcal{X}$ and consider the half-line from $\mathbf{t}(\mathbf{x}_1)$ through \mathbf{t}_0 . If this half-line crosses $\mathbf{t}(\mathcal{X})$ at some other point $\mathbf{t}(\mathbf{x}_2)$ after crossing \mathbf{t}_0 , we can write \mathbf{t}_0 as a convex combination of $\mathbf{t}(\mathbf{x}_1)$ and $\mathbf{t}(\mathbf{x}_2)$, which implies that $\langle \mathbf{u}, \mathbf{t}(\mathbf{x}_1) \rangle = \langle \mathbf{u}, \mathbf{t}(\mathbf{x}_2) \rangle = \langle \mathbf{u}, \mathbf{t}_0 \rangle$. If this half-line does not contain any other point in $\mathbf{t}(\mathcal{X})$, this whole half-line must be contained in $\text{int}(\mathcal{K})$. Take any other point $\mathbf{t}(\mathbf{x}_2) \in \text{int}(\mathcal{K})$ on the half-line after crossing \mathbf{t}_0 . By definition $\mathbf{t}(\mathbf{x}_2)$ is a convex combination of

some finitely many points in $\mathbf{t}(\mathcal{X})$, which allows us to write \mathbf{t}_0 as a convex combination of $\mathbf{t}(\mathbf{x}_1)$ and some other points in $\mathbf{t}(\mathcal{X})$. We again conclude that $\langle \mathbf{u}, \mathbf{t}(\mathbf{x}_1) \rangle = \langle \mathbf{u}, \mathbf{t}_0 \rangle$. In this way we showed that for an arbitrary $\mathbf{x} \in \mathcal{X}$, $\langle \mathbf{u}, \mathbf{t}(\mathbf{x}) \rangle = \langle \mathbf{u}, \mathbf{t}_0 \rangle$, which contradicts the minimality of our exponential family. \square

Proposition 1.3.6 shows that $\boldsymbol{\mu}$ can be used as an alternative parametrization of the exponential family. For example, in the mean zero Gaussian distribution the mean parametrization is given by the covariance matrix Σ (or more precisely by $-\frac{1}{2}\Sigma$). In the next section, we discuss a whole range of suitable parametrizations.

1.4 Marginal and conditional distributions*

We will consider partitioning of the sufficient statistics \mathbf{t} into \mathbf{u} and \mathbf{v} , $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ with the corresponding partition of $\boldsymbol{\theta} = (\boldsymbol{\theta}_u, \boldsymbol{\theta}_v)$ and $\boldsymbol{\mu} = (\boldsymbol{\mu}_u, \boldsymbol{\mu}_v)$. We consider two basic examples.

Example 1.4.1. We have shown that $X \sim N(\mu, \sigma)$ forms an exponential family with $\mathbf{t}(x) = (x, -x^2/2)$, $\boldsymbol{\theta} = (\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2})$ and $\boldsymbol{\mu} = (\mu, -\frac{1}{2}(\mu^2 + \sigma^2))$. Given a sample x_i for $i = 1, \dots, n$, we get an exponential family with the same canonical parameter and the sufficient statistics $\mathbf{t}(\mathbf{x}_{1:n}) = (\sum_i x_i, -\frac{1}{2}\sum_i x_i^2)$ and the mean parameter $(n\mu, -\frac{n}{2}(\mu^2 + \sigma^2))$. Here we could take $\mathbf{u}(x) = \sum_i x_i$ and $\mathbf{v}(x) = -\frac{1}{2}\sum_i x_i^2$ or the other way around.

Example 1.4.2. In the multivariate Gaussian case we have shown that $\mathbf{t}(\mathbf{x}) = -\frac{1}{2}\mathbf{x}\mathbf{x}^\top$, $\boldsymbol{\theta} = \mathbf{K}$, $\boldsymbol{\mu} = -\frac{1}{2}\Sigma$. Fix any subset $E \subset \{(i, i) : i = 1, \dots, m\} \cup \{(i, j) : 1 \leq i < j \leq m\}$. This corresponds to fixing some entries of a symmetric matrix. We could take $\mathbf{u} = -\frac{1}{2}(x_i x_j)_{ij \in E}$ and $\mathbf{v} = -\frac{1}{2}(x_i x_j)_{ij \notin E}$.

Recall the formula for the distribution of $\mathbf{t}(\mathbf{X})$ as given in Proposition 1.2.7.

Proposition 1.4.3 (Marginal distribution). *In a regular exponential family with $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ and $\boldsymbol{\theta} = (\boldsymbol{\theta}_u, \boldsymbol{\theta}_v)$ the marginal model for \mathbf{u} is a regular exponential family for each given $\boldsymbol{\theta}_v$, depending on $\boldsymbol{\theta}_v$ but with the same parameter space for its mean value parameter $\boldsymbol{\mu}_u$.*

Proof. The marginal distribution for \mathbf{u} is obtained by integrating \mathbf{v} out:

$$\begin{aligned} f(\mathbf{u}; \boldsymbol{\theta}) &= \int g(\mathbf{u}, \mathbf{v}) \exp \left\{ \langle \boldsymbol{\theta}_u, \mathbf{u} \rangle + \langle \boldsymbol{\theta}_v, \mathbf{v} \rangle - A(\boldsymbol{\theta}_u, \boldsymbol{\theta}_v) \right\} d\mathbf{v} \\ &= \exp \left\{ \langle \boldsymbol{\theta}_u, \mathbf{u} \rangle - A(\boldsymbol{\theta}_u, \boldsymbol{\theta}_v) \right\} \left(\int g(\mathbf{u}, \mathbf{v}) \exp \{ \langle \boldsymbol{\theta}_v, \mathbf{v} \rangle \} d\mathbf{v} \right). \end{aligned}$$

For any fixed $\boldsymbol{\theta}_v$ this has the form of a regular exponential family. This exponential family has canonical parameter $\boldsymbol{\theta}_u$ but the space of

canonical parameters will typically depend on θ_v (it is an intersection of Θ with $\theta_v = \text{fixed}$). However, by Proposition 1.3.6, the mean parameter space is always the same and equal to the interior of the convex hull of $\mathbf{u}(\mathcal{X})$, which is equal to the projection of $\mu(\Theta)$ on the \mathbf{u} coordinates (i.e. μ_u). \square

Proposition 1.4.4 (Conditional distribution). *With the same setting as in Proposition 1.4.3, the conditional model for \mathbf{x} given \mathbf{u} (and thus also for \mathbf{v} given \mathbf{u}) is a regular exponential family with canonical statistics \mathbf{v} . The conditional model depends on \mathbf{u} but with one and the same canonical parameter θ_v as in the joint model.*

Proof. We have

$$f(\mathbf{x}|\mathbf{u}; \boldsymbol{\theta}) = \frac{f(\mathbf{u}, \mathbf{x}; \boldsymbol{\theta})}{f(\mathbf{u}; \boldsymbol{\theta})} = \frac{h(\mathbf{x}) \exp\{\langle \boldsymbol{\theta}_v, \mathbf{v}(\mathbf{x}) \rangle\}}{\int g(\mathbf{u}, \mathbf{v}) \exp\{\langle \boldsymbol{\theta}_v, \mathbf{v}(\mathbf{x}) \rangle\} d\mathbf{v}}.$$

For the fixed value of \mathbf{u} , the expression in the denominator does not depend on \mathbf{x} but only on $\boldsymbol{\theta}_v$ and so it represents the normalizing constant of this distribution. Note $f(\mathbf{x}|\mathbf{u}; \boldsymbol{\theta})$ is defined only for those \mathbf{x} for which $\mathbf{t}(\mathbf{x}) = (\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x})) \in \mathbf{t}(\mathcal{X})$, and thus the space of sufficient statistics depends on \mathbf{u} . However, the canonical parameter is the same, $\boldsymbol{\theta}_v$ in the projection Θ_v of Θ on the coordinates $\boldsymbol{\theta}_v$. To get $f(\mathbf{v}|\mathbf{u}; \boldsymbol{\theta})$ we only substitute $h(\mathbf{x})$ for $g(\mathbf{u}, \mathbf{v})$ above but otherwise the argument is the same. \square

Explicit calculations with these marginal and conditional distributions are typically hard but the two results above are important in guiding our analysis.

Example 1.4.5. *Consider the univariate Gaussian example discussed in Example 1.4.1 and the induced distribution of the sample $x_{1:n}$. Let $u = \sum_i x_i$ and $v = -\frac{1}{2} \sum_i x_i^2$ and recall that the distribution has canonical parameters $(\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2})$. The marginal distribution of $\sum_i x_i^2$ is proportional to a noncentral χ^2 and does not in general form an exponential family unless we fix the value of $\frac{\mu}{\sigma^2}$. Next, consider instead its conditional distribution given $\sum_i x_i = n\bar{x}$. This may appear quite complicated but Proposition 1.4.4 suggests that it may be still tractable. Given \bar{x} , $\sum_i x_i^2$ differs only by an additive constant from $\sum x_i^2 - n\bar{x}^2 = (n-1)s^2$, where we used the standard notation*

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (1.13)$$

Thus, it is enough to characterize the distribution of $(n-1)s^2$ given \bar{x} . It is well-known that s^2 is independent of \bar{x} , and that the distribution of $(n-1)s^2$ is proportional (by σ^2) to a (central) χ_{n-1}^2 . From the explicit form of a χ^2 it is easily seen that the conditional distribution forms an exponential family.

← Exercise 1.9.8

As a corollary of the above two results we obtain a useful result on a range of possible alternative parametrizations. For a given split $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ consider the vector $(\boldsymbol{\mu}_u, \boldsymbol{\theta}_v)$ with $\boldsymbol{\mu}_u \in \mu_u(\Theta)$ and $\boldsymbol{\theta}_v \in \Theta_v$. Here by Θ_v we denote the projection of Θ on the coordinates $\boldsymbol{\theta}_v$ and by $\mu_u(\Theta)$ we mean the projection of $\mu(\Theta)$ on the coordinates $\boldsymbol{\mu}_u$.

Proposition 1.4.6. *The mixed parametrization $(\boldsymbol{\mu}_u, \boldsymbol{\theta}_v)$ is a valid parametrization with the parameter space $\mu_u(\Theta) \times \Theta_v$ (variational independence!). The Fisher information for $(\boldsymbol{\mu}_u, \boldsymbol{\theta}_v)$ is*

$$I(\boldsymbol{\mu}_u, \boldsymbol{\theta}_v) = \begin{bmatrix} (\Sigma_{uu})^{-1} & \mathbf{0} \\ \mathbf{0} & ((\Sigma^{-1})_{vv})^{-1} \end{bmatrix},$$

where $\Sigma = \text{var}(\mathbf{t})$ and $\Sigma_{uu} = \text{var}(\mathbf{u})$. The same formula holds for the observed information in the MLE, $J(\bar{\boldsymbol{\mu}}_u, \bar{\boldsymbol{\theta}}_v)$.

Proof. Fix an exponential family with canonical statistics $\mathbf{t}(\mathbf{x})$ and canonical parameter $\boldsymbol{\theta} \in \Theta$. By Proposition 1.2.7, the distribution of $\mathbf{t}(\mathbf{X})$ is an exponential family with the same canonical parameter. This distribution is uniquely defined by the marginal distribution of \mathbf{u} and the conditional distribution of \mathbf{v} given \mathbf{u} . The latter forms an exponential family with canonical parameter $\boldsymbol{\theta}_v \in \Theta_v$ by Proposition 1.4.4. Now fix $\boldsymbol{\theta}_v$, which corresponds to fixing the conditional distribution of \mathbf{v} given \mathbf{u} . By Proposition 1.4.3, the marginal distribution of \mathbf{u} is an exponential family with the mean parameter $\boldsymbol{\mu}_u$. By Proposition 1.3.6 the range of this mean parameter is the interior of the convex hull of $\mathbf{u}(\mathcal{X})$ (independent on $\boldsymbol{\theta}_v$). This is precisely the projection of $\mu(\Theta)$ on $\boldsymbol{\mu}_u$ and this shows that the map $\boldsymbol{\theta} \mapsto (\boldsymbol{\mu}_u, \boldsymbol{\theta}_v)$ is one-to-one with range $\mu_u(\Theta) \times \Theta_v$. For the proof of the second statement see ⁹. □

Example 1.4.7. *Consider the multivariate Gaussian distribution in Example 1.1.6 with $m = 2$. We have $K = \Sigma^{-1}$, that is,*

$$K = \begin{bmatrix} K_{11} & K_{12} \\ K_{12} & K_{22} \end{bmatrix} = \frac{1}{\Sigma_{11}\Sigma_{22} - \Sigma_{12}^2} \begin{bmatrix} \Sigma_{22} & -\Sigma_{12} \\ -\Sigma_{12} & \Sigma_{11} \end{bmatrix}.$$

The canonical parameters are (K_{11}, K_{22}, K_{12}) and the mean parameters are $-\frac{1}{2}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12})$. The constraints on the canonical parameters are $K_{11} > 0, K_{22} > 0, K_{11}K_{22} > K_{12}^2$ (namely, K is positive definite). The constraints on the mean parameter follow from the constraints on Σ : $\Sigma_{11} > 0, \Sigma_{22} > 0, \Sigma_{11}\Sigma_{22} > \Sigma_{12}^2$ (Σ is positive definite). Consider a mixed parametrization $(-\frac{1}{2}\Sigma_{11}, -\frac{1}{2}\Sigma_{22}, K_{12})$. The projection $\mu_u(\Theta)$ is simply $(-\infty, 0)^2$. The projection, Θ_v is the whole real line \mathbb{R} (irrespective of the value of K_{12} we can set K_{11}, K_{22} big enough for K to be positive definite).

⁹ Rolf Sundberg. *Statistical modelling by exponential families*, volume 12 of *Institute of Mathematical Statistics Textbooks*. Cambridge University Press, Cambridge, 2019

← Exercise 1.9.9

← Exercise 1.9.10

By Proposition 1.4.6, for every choice of $\Sigma_{11} > 0, \Sigma_{22} > 0$ and $K_{12} \in \mathbb{R}$, there will be a unique positive definite matrix Σ with the prescribed diagonal entries Σ_{11}, Σ_{22} and such that $-\frac{1}{\Sigma_{11}\Sigma_{22}-\Sigma_{12}^2}\Sigma_{12} = K_{12}$. In this case, this can be checked directly. To compute the corresponding Σ_{12} , we need to solve the quadratic equation

$$K_{12}\Sigma_{12}^2 - \Sigma_{12} - K_{12}\Sigma_{11}\Sigma_{22} = 0.$$

There are two real solutions but only one of them, namely,

$$\Sigma_{12} = \frac{1 - \sqrt{1 + 2K_{12}^2\Sigma_{11}\Sigma_{22}}}{2K_{12}}$$

results in a positive definite Σ .

1.5 Conditional inference for canonical parameter*

Suppose ψ is the parameter of interest, where (λ, ψ) is a transformation of θ . For simplicity we focus on the case when $\psi = \theta_v$ and $\lambda = \mu_u$ is regarded as nuisance parameter. As shown in Proposition 1.4.6, $\lambda = \mu_u = \mathbb{E}_\theta(\mathbf{u})$ is the preferable nuisance parameter (rather than θ_u), since θ_v and μ_u are variation independent and information orthogonal.

Proposition 1.5.1 (Conditionality principle for full families). *Statistical inference about the canonical parameter component θ_v in presence of the nuisance parameter $\lambda = \mu_u = \mathbb{E}_\theta(\mathbf{u})$ should be made conditional on \mathbf{u} , that is, the conditional model for \mathbf{x} or \mathbf{v} given \mathbf{u} .*

This is only a recommendation so rather than providing a formal proof we motivate this statement informally.

Motivation. The likelihood for (μ_u, θ_v) factorizes as

$$L(\mu_u, \theta_v; \mathbf{t}) = L_1(\mu_u, \theta_v; \mathbf{u})L_2(\theta_v; \mathbf{v}|\mathbf{u}) \quad (1.14)$$

where the two parameters are variation independent. In some cases L_1 depends only on μ_u , in which case

$$L(\mu_u, \theta_v; \mathbf{t}) = L_1(\mu_u; \mathbf{u})L_2(\theta_v; \mathbf{v}|\mathbf{u}).$$

Then it is clear that there is no information about θ_v in the first factor L_1 and the argument for the principle is compelling.

However, even when L_1 depends on θ_v , there is really no information about θ_v in \mathbf{u} . Indeed, note first that \mathbf{u} and μ_u have the same dimension, and that \mathbf{u} serves as an estimator (the MLE) of μ_u , whatever is the value of θ_v . This means that the information in \mathbf{u} about (μ_u, θ_v)

is totally consumed in the estimation of μ_u . Furthermore, the estimated value of μ_u does not provide any information about θ_v , and μ_u would not do so even if it were known, due to the variation independence between μ_u and θ_v . Thus, the first factor L_1 contributes only information about μ_u . \square

Example 1.5.2 (Conditional independence for a Gaussian sample).

Suppose we want to make inference about σ^2 , or σ . Then we are led to consider the conditional distribution of $\sum_i x_i^2$, given \bar{x} , that depends on σ alone (c.f. Example 1.4.5 and Proposition 1.4.4). The marginal distribution of \bar{x} depends on both μ and σ^2 , so the joint and conditional likelihoods are different functions of σ^2 .

As we have seen in Example 1.4.5, in the conditional approach $n\bar{x}^2$ is a constant, and after subtraction of this empirical constant from $\sum_i x_i^2$ we are led to the use of the statistic s^2 . Now, we already know that \bar{x} and s^2 are independent, so the even simpler result is that the inference should be based on the marginal model for $(n-1)s^2/\sigma^2$, with its χ_{n-1}^2 -distribution. In particular this leads to the conditional and marginal ML estimator $\hat{\sigma}^2 = s^2$, which differs by the factor $n/(n-1)$ from the MLE in the joint model (with denominator n).

1.6 Kullback-Leibler divergence

The Fenchel conjugate of the cumulant function A is the function

$$A^*(\mathbf{t}) = \sup\{\langle \boldsymbol{\theta}, \mathbf{t} \rangle - A(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathbb{R}^d\}.$$

For regular exponential families $A^*(\mathbf{t}) < \infty$ if and only if $\mathbf{t} \in \mu(\Theta)$. The function A^* is convex as a supremum of linear functions¹⁰ and, in fact, strictly convex. If $\mathbf{t} \in \mu(\Theta)$, the unique optimizer of the log-likelihood is $\theta(\mathbf{t})$, where $\theta : M \rightarrow \Theta$ is the inverse of the map $\mu : \Theta \rightarrow M$; see Proposition 1.3.6. It follows that, for $\mathbf{t} \in \mu(\Theta)$,

$$A^*(\mathbf{t}) = \langle \theta(\mathbf{t}), \mathbf{t} \rangle - A(\theta(\mathbf{t})) \quad (1.15)$$

or alternatively, for $\boldsymbol{\theta} \in \Theta$,

$$A^*(\mu(\boldsymbol{\theta})) = \langle \boldsymbol{\theta}, \mu(\boldsymbol{\theta}) \rangle - A(\boldsymbol{\theta}) \quad (1.16)$$

implying in particular that A^* is smooth, since μ and A are both smooth. By composite differentiation in (1.15), since $\mu(\theta(\mathbf{t})) = \mathbf{t}$, we obtain

$$\nabla A^*(\mathbf{t}) = \theta(\mathbf{t}) + \nabla \theta(\mathbf{t}) \cdot \mathbf{t} - \nabla \theta(\mathbf{t}) \cdot \mu(\theta(\mathbf{t})) = \theta(\mathbf{t}),$$

where the gradient is taken with respect to \mathbf{t} .

¹⁰ see Proposition B.2.5.

For two distributions over some state-space \mathcal{X} with densities p, q with respect to some measure μ , the **Kullback-Leibler divergence** $K(p, q)$ is defined as

$$K(p, q) := \int_{\mathcal{X}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \mu(d\mathbf{x}).$$

The following result is well-known and it is a direct application of the Jensen's inequality stated formally in Theorem B.2.3.

Proposition 1.6.1. *We have $K(p, q) \geq 0$ with equality if and only if $p = q$ almost surely.*

Proof. We use the fact that $-\log y$ is a strictly convex function. By Theorem B.2.3,

$$0 = -\log \mathbb{E}_p \frac{q(X)}{p(X)} \leq -\mathbb{E}_p \log \frac{q(X)}{p(X)} = \mathbb{E}_p \log \frac{p(X)}{q(X)} = K(p, q)$$

with equality if and only if $q(X)/p(X)$ is constant almost surely. Since p, q are both densities, this is possible if and only if they are equal almost surely. \square

Given two distributions $\mathbb{P}_{\theta_1}, \mathbb{P}_{\theta_2}$ in the given exponential family we write the corresponding Kullback-Leibler divergence as $K(\theta_2, \theta_1)$. We easily check that

$$K(\theta_1, \theta_2) = A(\theta_2) - A(\theta_1) - \langle \theta_2 - \theta_1, \mu_1 \rangle. \quad (1.17)$$

Proposition 1.6.2. *Consider two distributions in the exponential family (1.1), one with the mean parameter $\mu_1 \in \mu(\Theta)$ and the other with canonical parameter $\theta_2 \in \Theta$. If this exponential family is regular, the Kullback–Leibler divergence between these two distributions is*

$$K(\mu_1, \theta_2) = -\langle \mu_1, \theta_2 \rangle + A^*(\mu_1) + A(\theta_2). \quad (1.18)$$

The Kullback–Leibler divergence is well defined and nonnegative over $\mu(\Theta) \times \Theta$. Moreover, $K(\mu_1, \theta_2) = 0$ if and only if $\mu_1 = \mu(\theta_2)$.

Proof. We leave it as an exercise. \square

Since $\mu_1 = \nabla A(\theta_1)$, (1.17) has another interpretation as the Bregman divergence (defined by the function A between θ_2 and θ_1). Look this up!

← Exercise 1.9.11

The reason to express the Kullback–Leibler distance in terms of μ_1 and θ_2 rather than θ_1, θ_2 (as usually done in the literature) is that we wish to exploit the following basic result.

Proposition 1.6.3. *The Kullback–Leibler divergence $K(\mu_1, \theta_2)$ is strictly convex both in μ_1 and in θ_2 .*

Proof. This follows directly from (1.18) and the fact that both $A(\theta)$ and $A^*(\mu)$ are strictly convex functions \square

This set-up has been exploited in various places. See, for example, Section 5 in ¹¹. Another important application is in situations when a statistical submodel is given by affine restrictions on the mean parameter. Note that the MLE can be equivalently defined as the minimizer of $K(\bar{\mu}_n, \theta)$. We have a parallel definition, when the dual MLE is given as the minimizer of $K(\mu, \theta(\bar{\mu}_n))$.

Example 1.6.4 (Behrens-Fisher problem). *The Behrens-Fisher problem is concerned with testing the difference between the means of two normally distributed populations when the variances of the two populations are not assumed to be equal, based on two independent samples. Since the hypothesis is linear in the mean parameter, this problem can be addressed with the dual MLE; see ¹² for details.*

We finish this chapter with one of the most fundamental results motivating exponential families. The maximum entropy principle states that under uncertainty, one should take a model which maximizes the entropy subject to constraints on the known features about the system. We show that the exponential family arises naturally if the constraint is given by the expected value of some statistics.

Recall that for a distribution \mathbb{P} that admits a density function $p(x)$ with respect to the base measure μ , the entropy $H_{\mathbb{P}}$ of \mathbb{P} is

$$H_{\mathbb{P}} = -\mathbb{E}_{\mathbb{P}} \log p(X) = -\int \log p(x)p(x)\mu(dx).$$

For notational simplicity, in what follows consider the exponential family (1.2), where the function $h(\mathbf{x})$ has been incorporated into the base measure μ . In this case each distribution in this family has the same support, which is equal to the support of μ . For such an exponential family \mathbb{P}_{θ} we have

$$H_{\mathbb{P}_{\theta}} = -\langle \theta, \mu(\theta) \rangle + A(\theta) = -A^*(\mu(\theta)), \quad (1.19)$$

where the second equality follows from (1.16).

Consider the problem of maximizing the entropy for *all* distributions \mathbb{P} that admit a density function absolutely continuous with respect to μ and $\mathbb{E}_{\mathbb{P}}(\mathbf{t}(X)) = \mathbf{t}_0$

$$\text{maximize } H_{\mathbb{P}} \quad \text{s.t. } \mathbb{P} \sim \mu, \quad \mathbb{E}_{\mathbb{P}} \mathbf{t}(X) = \mathbf{t}_0. \quad (1.20)$$

The following result provides an important characterization of exponential families.

Theorem 1.6.5. *Consider the exponential family (1.2). Suppose there exists $\theta_0 \in \Theta$ such that $\nabla A(\theta_0) = \mathbf{t}_0$. Then for any distribution \mathbb{P} , which satisfies condition of (1.20), we have*

$$H_{\mathbb{P}_{\theta_0}} - H_{\mathbb{P}} = K(\mathbb{P}, \mathbb{P}_{\theta_0}).$$

Thus, \mathbb{P}_{θ_0} is the unique solution to (1.20).

¹¹ Martin J Wainwright and Michael I. Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008

¹² E Susanne Christensen. Statistical properties of I-projections within exponential families. *Scandinavian Journal of Statistics*, pages 307–318, 1989

Proof. Note that $\nabla A(\boldsymbol{\theta}_0) = \mathbf{t}_0$ is equivalent to $\mu(\boldsymbol{\theta}_0) = \mathbb{E}_{\boldsymbol{\theta}_0} \mathbf{t}(X) = \mathbf{t}_0$. Let $p(\mathbf{x})$ be the density of \mathbb{P} with respect to μ . Consider

$$\begin{aligned} K(\mathbb{P}, \mathbb{P}_{\boldsymbol{\theta}_0}) &= \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{f(\mathbf{x}; \boldsymbol{\theta}_0)} \mu(d\mathbf{x}) \\ &= -H_{\mathbb{P}} - \int (\langle \boldsymbol{\theta}_0, \mathbf{t}(\mathbf{x}) \rangle - A(\boldsymbol{\theta}_0)) p(\mathbf{x}) \mu(d\mathbf{x}) \\ &= -H_{\mathbb{P}} - \langle \boldsymbol{\theta}_0, \mathbf{t}_0 \rangle + A(\boldsymbol{\theta}_0) \\ &\stackrel{(1.19)}{=} H_{\mathbb{P}_{\boldsymbol{\theta}_0}} - H_{\mathbb{P}}. \end{aligned}$$

By Proposition 1.6.1, this shows that $H_{\mathbb{P}} \leq H_{\mathbb{P}_{\boldsymbol{\theta}_0}}$ with equality if and only if $\mathbb{P} = \mathbb{P}_{\boldsymbol{\theta}_0}$. This concludes the proof. \square

Example 1.6.6. Consider all distributions with the support \mathbb{R}^m and with the property that $\mathbb{E}X = \mu$, $\mathbb{E}(XX^\top) = \Sigma + \mu\mu^\top$. Among all such distributions, the multivariate normal $N(\mu, \Sigma)$ is the one that maximizes the entropy.

1.7 Generalized Linear Models

The generalized linear models are formulated based on the construction of exponential families. Here we provide only a basic treatment that explains the origin of the construction and the most important examples.

Consider the pairs $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$, where the input $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ are considered fixed and the outputs y_1, \dots, y_n are independent observations each from density

$$f(y_i; \mathbf{x}_i, \boldsymbol{w}, \sigma^2) = h(y_i, \sigma^2) \exp \left\{ \frac{1}{\sigma^2} (y_i \theta_i(\boldsymbol{w}) - A(\theta_i(\boldsymbol{w}))) \right\}.$$

where $\theta_i(\boldsymbol{w}) = \mathbf{x}_i^\top \boldsymbol{w}$ and σ^2 is called the dispersion term. From the standard theory, we get immediately that

$$\mu(\theta_i) := \mathbb{E}[Y_i | \mathbf{x}_i, \boldsymbol{w}, \sigma^2] = A'(\theta_i).$$

Indeed, the standard result show that for a fixed σ^2 , $\mathbb{E}(\frac{1}{\sigma^2} Y) = \frac{1}{\sigma^2} A'(\boldsymbol{\theta})$. In the same way, we argue that

$$V(\boldsymbol{\theta}, \sigma^2) := \text{var}_{\boldsymbol{\theta}}(Y_i | \mathbf{x}_i, \boldsymbol{w}, \sigma^2) = \sigma^2 A''(\boldsymbol{\theta}_i).$$

From now on we fix σ^2 and with no loss of generality we take $\sigma^2 = 1$. The log-likelihood function is

$$\ell_n(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^n \log f(y_i; \mathbf{x}_i, \boldsymbol{w}) = \left(\frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i \right)^\top \boldsymbol{w} - \frac{1}{n} \sum_{i=1}^n A(\mathbf{x}_i^\top \boldsymbol{w}) + \text{const.}$$

Since a composition of a convex and a linear function is convex, we conclude that $\ell_n(\boldsymbol{w})$ is a concave function.

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the data matrix with rows \mathbf{x}_i and let \mathbf{y} be the vectors with entries y_i . It can be shown that the unique optimizer of $\ell_n(\mathbf{w})$ (if it exists), must satisfy the likelihood equations

← Exercise 1.9.13

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top A'(\mathbf{X}\hat{\mathbf{w}}), \quad (1.21)$$

where the function A' is applied elementwise to the vector $\mathbf{X}\hat{\mathbf{w}}$. In machine learning it is customary to call A' in this context an activation function.

We now discuss a bunch of basic examples.

Example 1.7.1 (Linear regression). Consider the univariate Gaussian distribution in Example 1.1.4. Suppose now that σ^2 is fixed and the model is parametrized only by the mean μ . We can rewrite the density as

$$f(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{y^2}{2\sigma^2}} \exp\left\{\frac{1}{\sigma^2}(y\mu - \frac{\mu^2}{2})\right\} = h(y, \sigma^2) \exp\left\{\frac{1}{\sigma^2}(y\mu - A(\mu))\right\},$$

where $h(y, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{y^2}{2\sigma^2}}$ and $A(\mu) = \mu^2/2$. If we model the mean μ as a linear function of a vector \mathbf{x} , $\mu = \mathbf{w}^\top \mathbf{x}$ we get the standard Gaussian linear regression

$$f(y|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y - \mathbf{w}^\top \mathbf{x})^2\right\}.$$

We easily see that this is a generalized linear model as defined above. As $A'(\mu) = \mu$ and $A''(\mu) = 1$, we get $\mathbb{E}(Y|\mathbf{x}, \mathbf{w}, \sigma^2) = \mu = \mathbf{w}^\top \mathbf{x}$, $\text{var}(Y|\mathbf{x}, \mathbf{w}, \sigma^2) = \sigma^2$ and, by (1.21), the MLE equations are $\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X}\mathbf{w}$, which are identical to the OLS estimating equations.

Example 1.7.2 (Binomial regression). The binomial distribution for the number of successes in n trials, $y \in \{0, \dots, n\}$ has the probability mass function

$$\text{Bin}(y; n, p) = \binom{n}{y} p^y (1-p)^{n-y} = \binom{n}{y} \exp\left\{y \log \frac{p}{1-p} + n \log(1-p)\right\}.$$

If $\theta = \log \frac{p}{1-p}$ then $p = \frac{e^\theta}{1+e^\theta}$, which is typically called a sigmoid function and denoted by $\sigma(\theta)$. Moreover, $A(\theta) = n \log(1+e^\theta)$ and so $A'(\theta) = n\sigma(\theta)$. Now we use this distribution for the GLM setup. If the response variable is the number of successes in n trials, $y \in \{0, \dots, n\}$, we can use binomial regression, which is defined by

$$f(y; \mathbf{x}, N, \mathbf{w}) = \text{Bin}(y; n, \sigma(\mathbf{w}^\top \mathbf{x})),$$

where the logistic regression becomes a special case with $n = 1$. This clearly has the right form and $\mathbb{E}(Y) = n\theta$, $\text{var}(Y) = np(1-p)$. By Exercise 1.9.13 the likelihood equations are

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \sigma(\mathbf{X}\hat{\mathbf{w}}).$$

Example 1.7.3 (Poisson regression). *The Poisson distribution is a distribution over $\mathcal{X} = \{0, 1, 2, \dots\}$ with the probability mass function*

$$f(y; \mu) = \frac{\mu^y}{y!} e^{-\mu} \quad y \in \mathcal{X}.$$

This is an exponential family with mean parameter μ ($\mathbb{E}Y = \mu$) and canonical parameter $\theta = \log(\mu)$. In Poisson regression we take $\theta = \mathbf{w}^\top \mathbf{x}$.

Modelling the canonical parameters as a linear function of external variables is not the only choice. For example, for the Bernoulli distribution in Example 1.7.2 gives the logistic regression. An alternative approach is to model $p = \Phi(\mathbf{w}^\top \mathbf{x})$, which gives the probit regression. Generalized linear models with non-canonical link functions correspond to curved exponential families.

1.8 Diaconis-Ylvisaker conjugate priors

An important advantage of exponential families over more general classes of distributions is that they admit explicit conjugate prior for Bayesian computations. The conjugate prior measure for the exponential family (1.1) is given by the density (w.r.t. the Lebesgue measure) of the form

$$\pi(\boldsymbol{\theta}) = C \exp\{\langle \boldsymbol{\theta}, \boldsymbol{\tau} \rangle - n_0 A(\boldsymbol{\theta})\}, \quad \boldsymbol{\tau} \in \mathbb{R}^d, n_0 \geq 0. \quad (1.22)$$

Note that $\pi(\boldsymbol{\theta}) \equiv 0$ outside of Θ because there $A(\boldsymbol{\theta}) \equiv +\infty$. It can be also shown that the distribution is normalizable if $n_0 > 0$ and $\boldsymbol{\tau}/n_0 \in \mathcal{K} = \text{conv}(T(\mathcal{X}))$ ¹³.

Proposition 1.8.1. *For a regular exponential family, consider the conjugate prior in (1.22). If $\boldsymbol{\mu}$ is the mean parameter then we have $\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\mu}] = \frac{\boldsymbol{\tau}}{n_0}$.*

Proof sketch. We use the fact that $\boldsymbol{\mu} = \nabla A(\boldsymbol{\theta})$. We have

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\tau} - n_0 \nabla A(\boldsymbol{\theta})] &= \int_{\Theta} (\boldsymbol{\tau} - n_0 \nabla A(\boldsymbol{\theta})) C \exp\{\langle \boldsymbol{\theta}, \boldsymbol{\tau} \rangle - n_0 A(\boldsymbol{\theta})\} d\boldsymbol{\theta} \\ &= \int_{\Theta} \nabla (C \exp\{\langle \boldsymbol{\theta}, \boldsymbol{\tau} \rangle - n_0 A(\boldsymbol{\theta})\}) d\boldsymbol{\theta} \\ &= \int_{\Theta} \nabla \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \end{aligned}$$

It can be shown that the latter integral is equal to zero. This result is a consequence of Green's theorem (a general form of the fundamental theorem of calculus). A rigorous proof of this result was first presented in Diaconis and Ylvisaker (1979); for a simplified proof see ¹⁴.

¹³ Persi Diaconis and Donald Ylvisaker. Conjugate priors for exponential families. *The Annals of statistics*, pages 269–281, 1979

¹⁴ Lawrence D. Brown. *Fundamentals of statistical exponential families with applications in statistical decision theory*, volume 9 of *Institute of Mathematical Statistics Lecture Notes—Monograph Series*. Institute of Mathematical Statistics, Hayward, CA, 1986

If the prior is of the form (1.22) then the posterior $\pi(\boldsymbol{\theta}|\mathbf{x}_{1:n})$ satisfies

$$\begin{aligned}\pi(\boldsymbol{\theta}|\mathbf{x}_{1:n}) &\propto \pi(\boldsymbol{\theta}) \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\theta}) \\ &= C \prod_{i=1}^n h(\mathbf{x}_i) \exp\left\{\langle \boldsymbol{\theta}, \tau + \sum_{i=1}^n \mathbf{t}(\mathbf{x}_i) \rangle - (n + n_0)A(\boldsymbol{\theta})\right\}.\end{aligned}$$

And so it has the same general form as the prior, where (τ, n_0) is replaced with $(\tau + n\bar{\mu}_n, n_0 + n)$. By Proposition 1.8.1, the Bayes estimator of the mean parameter is

$$\mathbb{E}[\boldsymbol{\mu}|\mathbf{x}_{1:n}] = \frac{\tau + n\bar{\mu}_n}{n_0 + n} = \lambda \frac{\tau}{n_0} + (1 - \lambda)\bar{\mu}_n,$$

where $\lambda = \frac{n_0}{n_0 + n}$ and so the Bayes estimator is a convex combination of the MLE $\bar{\mu}_n$ and the prior mean τ/n_0 .

We can easily define conjugate priors over non-canonical parameter. For example, to get a conjugate prior over the mean parameter we simply change $\boldsymbol{\theta}$ with $\boldsymbol{\theta}(\boldsymbol{\mu})$ in (1.22) and change the definition of C so that the corresponding expression integrates to 1. Note that this is not the same as the density obtained through the change of variable formula.

Example 1.8.2. *If we do it for the Bernoulli distribution, the conjugate prior for the mean parameter is the Beta distribution.*

1.9 Exercises

Exercise 1.9.1. *Prove formally, using (1.4), that the space of canonical parameters in the centered multivariate Gaussian distribution is S_+^m .*

Exercise 1.9.2. *Consider a multivariate Gaussian distribution with general mean vector $\boldsymbol{\mu} \in \mathbb{R}^m$. Show that the canonical parameter space is $\mathbb{R}^m \times S_+^m$ with canonical parameters $(K\boldsymbol{\mu}, K)$ and that the sufficient statistics is $(\mathbf{x}, -\frac{1}{2}\mathbf{x}\mathbf{x}^\top)$.*

Exercise 1.9.3. *Show that the distribution of $\mathbf{x}_{1:n}$ in Proposition 1.1.7 is of exponential type with the same canonical space Θ . What is the sufficient statistics? (c.f. Example 1.1.4)*

Exercise 1.9.4 (Ising model on a bipartite graph). *Consider the Ising model on the bipartite graph G with m nodes X_1, \dots, X_m and n nodes Y_1, \dots, Y_n such that G has mn edges connecting each X_i with each Y_j . Show that to compute the conditional distribution of $Y = (Y_1, \dots, Y_n)$ given $X = (X_1, \dots, X_m)$ we essentially need to: (i) apply a linear function to X , (ii) apply an activation function to each element of the resulting vector. (Does it ring a bell?)*

Exercise 1.9.5 (Gaussian graphical models). Let G be a graph over m nodes representing m random variables $X_i \in \mathbb{R}$ for $i = 1, \dots, m$ with mean zero. Describe the distribution satisfying

$$f(\mathbf{x}; \boldsymbol{\theta}) \propto \exp\left\{\sum_{i=1}^m \theta_{ii} x_i^2 + \sum_{ij \in G} \theta_{ij} x_i x_j\right\}.$$

Exercise 1.9.6. Consider an exponential family with sufficient statistics $\mathbf{t}(\mathbf{x})$ and canonical parameter $\boldsymbol{\theta} \in \Theta$. Consider now a model whose parameter space is $\mathcal{L} \cap \Theta$ for some affine subspace $\mathcal{L} \subseteq \mathbb{R}^d$, such that $\mathcal{L} \cap \Theta \neq \emptyset$. Show that it forms an exponential family with $\Theta' = \Theta \cap \mathcal{L}$ and sufficient statistics which is a linear transformation of $\mathbf{t}(\mathbf{x})$.

Exercise 1.9.7. Suppose the state-space \mathcal{X} is finite. Show that every exponential family over \mathcal{X} is regular.

Exercise 1.9.8. Suppose that $\mathbf{X} = (X_1, \dots, X_n)$ has i.i.d. components that are $N(\mu, \sigma^2)$. Show that the vector $\mathbf{X} - \bar{X}\mathbf{1}$ is independent of $\bar{X} = \frac{1}{n}\mathbf{1}^\top \mathbf{X}$. Use this fact to conclude that s^2 defined in (1.13) is independent of \bar{x} . Hint: This result follows from basic matrix algebra. Let $A = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$. Note that $A\mathbf{X} = \mathbf{X} - \bar{X}\mathbf{1}$ and that $\text{cov}(\mathbf{X}) = \sigma^2 I_n$. Finally, note that $(n-1)s^2 = \text{tr}(A\mathbf{x}\mathbf{x}^\top A^\top)$.

Exercise 1.9.9. Let E be any set of pairs of elements of $\{1, \dots, m\}$. Use Proposition 1.4.6 to show that for any two $A, B \in \mathbb{S}_+^m$ there exists a unique $X \in \mathbb{S}_+^m$ such that $X_{ij} = A_{ij}$ for $ij \in E$ and $(X^{-1})_{ij} = B_{ij}$ for $ij \notin E$.

Exercise 1.9.10. Consider the Gaussian distribution $N(\mu, \sigma^2)$. From the first principles, provide the two mixed parametrizations. Discuss their set of parameters and the corresponding Fisher information matrices.

Exercise 1.9.11. Prove Proposition 1.6.2.

Exercise 1.9.12. In the zero-mean Gaussian distribution $N(0, \Sigma)$ with $K = \Sigma^{-1}$, the log-likelihood function is

$$\ell(K; S_n) = \log \det K - \langle K, S_n \rangle,$$

where $S_n = \frac{1}{n} \sum_i x_i x_i^\top$ is the sample covariance matrix. Following Section 1.6 of the notes show that the dual log-likelihood, up to some additive constants, is

$$\check{\ell}(\Sigma; S_n) = \log \det \Sigma - \langle \Sigma, S_n^{-1} \rangle.$$

Consider the bivariate Gaussian distribution with mean zero and covariance

$$\Sigma = \begin{bmatrix} a & b \\ b & a \end{bmatrix}.$$

Compare in simulations the maximum likelihood estimator (\hat{a}, \hat{b}) of (a, b) with the dual maximum likelihood estimator (\check{a}, \check{b}) . Based on your simulations, what is the asymptotic behaviour of $\sqrt{n}(\hat{a} - \check{a})$?

Exercise 1.9.13. Consider a generalized linear model with canonical link function. Argue that the maximum likelihood estimator of the parameter \boldsymbol{w} , for data (y_i, \mathbf{x}_i) for $i = 1, \dots, n$, leads to a convex optimization problem. Show that the MLE satisfies

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top A'(\mathbf{X}\boldsymbol{w}),$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the data matrix, $\mathbf{y} = (y_1, \dots, y_n)$, and in $A'(\mathbf{X}\boldsymbol{w})$ the function A' is applied elementwise to the vector $\mathbf{X}\boldsymbol{w}$.

2

Statistical Decision Theory (1 week)

Statistical decision theory was covered in the first semester. Here I recall some of this material but in class we focus entirely on admissibility.

2.1 Decision theoretic framework*

Statistical decision theory, developed by Abraham Wald, Jerzy Neyman and others during the mid-20th century, provides an abstraction that allows for comparison of statistical procedures. Our decision theoretic framework is made up from the following ingredients:

• **A family of probabilistic models with parameterization** $\theta \in \Theta$.

We can think of this as a mapping from the parameter space Θ to a family of probability distributions $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$. We assume that all distributions in \mathcal{P} are absolutely continuous with respect to some underlying measure μ . Examples of parameter spaces and models include the following:

- $\Theta \subseteq \mathbb{R}^d$: exponential, Gaussian, other parametric models.
- $\Theta \subseteq$ (a function space): in nonparametric settings, e.g. the set of all twice-differentiable functions.

• **A decision procedure** In general this is a recipe that defines what action to take, given a set of observations (\mathcal{X}, μ) . The set of possible actions is denoted \mathcal{A} . Examples of possible actions might include:

- Accepting or rejecting a null hypothesis, $\mathcal{A} = \{0, 1\}$.
- Estimating a value for some model parameter θ , $\mathcal{A} = \Theta$.
- Selecting one family of models as “superior” to other models (model selection).

Formally, a (non-randomized) **decision rule** is a measurable function $\delta : \mathcal{X} \rightarrow \mathcal{A}$ or $\delta : \mathcal{X}^n \rightarrow \mathcal{A}$. An estimator δ is a particular kind

of decision procedure by which we estimate the value of θ from the observations. In this case $\mathcal{A} = \Theta$. Denote by \mathcal{D}_0 the set of all (non-randomized) decision rules.

• **A loss function**, which tells us how to evaluate different decision procedures. This is an extended mapping $L : \Theta \times \mathcal{A} \rightarrow [0, \infty]$, in which $L(\theta, a)$ represents the loss incurred by deciding a when θ is “true”. In many applications $L(\theta, a)$ is convex both in θ and in a (or jointly) but there are important examples where there is no convexity. Canonical examples of the loss functions are the convex loss functions $L(\theta, a) = \|\theta - a\|^2$, $L(\theta, a) = \|\theta - a\|$, and the 0/1-loss $L(\theta, a) = \mathbb{1}(\theta \neq a)$. Depending on the particular application we also consider hinge loss, Kullback-Leibler divergence, and many others.

We assume that $L(\theta, a)$ lower semicontinuous in a . This means that for every $\theta \in \Theta$ and every $t \in \mathbb{R}$ the set $\{a \in \mathcal{A} : L(\theta, a) \leq t\}$ is closed. By Exercise 2.7.1, such $L(\theta, a)$ is then also measurable in a^1 . In particular, if X has distribution \mathbb{P}_θ then $L(\theta, \delta(X))$ is \mathbb{P}_θ measurable.

We define the **risk function** as:

$$R(\theta, \delta) = \mathbb{E}_\theta L(\theta, \delta(X)).$$

When the sample space is continuous, a quadratic loss function is frequently used. In this case, the risk is simply the mean squared error.

Example 2.1.1. *Suppose a coin is being tossed and you are interested in the probability of getting heads. We can model this using a family of Bernoulli distributions on the binary outcome space: $X_i \sim \text{Ber}(\theta)$ with $\theta \in \Theta = (0, 1)$. Suppose we observe n replications with i.i.d results $X = (X_1, \dots, X_n) \in \{0, 1\}^n$ from \mathbb{P}_θ . We consider various rules $\delta_i(X)$ for estimating θ . In this case $\mathcal{A} = [0, 1]$. We use the quadratic loss $L(\theta, \delta(X)) = (\theta - \delta(X))^2$. and evaluate the associated risk function for each estimator:*

- The first estimator is the sample mean $\delta_1(X) = \frac{1}{n} \sum_{i=1}^n X_i$ with risk

$$R(\theta, \delta_1) = \frac{1}{n} \theta(1 - \theta).$$

- The second estimator is given by the constant value $\delta_2(X) = \frac{1}{2}$ with risk

$$R(\theta, \delta_2) = (\theta - \frac{1}{2})^2.$$

- The third estimator

$$\delta_3(X) = \frac{\sum_{i=1}^n X_i + 3}{n + 6}$$

with risk

$$R(\theta, \delta_3) = \frac{9 + (n - 36)\theta - (n - 36)\theta^2}{(n + 6)^2}.$$

← Exercise 2.7.1

¹ We think about $\mathcal{A} \subseteq \mathbb{R}^d$ as a measurable space with the underlying σ -field of Borel sets (generated by the open subsets). A function $f : \mathcal{A} \rightarrow \mathbb{R}$ is then measurable if the sets $\{a : f(a) \leq t\}$ are measurable.

The risks functions for the first and the last estimator are plotted in Figure 2.1. Intuitively, it seems clear that if we strongly believe that θ is close to 0.5, δ_3 outperforms δ_1 . Later we learn how to formalize that.

Example 2.1.2 (Benefits of bias). Suppose $X_i \sim U(0, \theta)$ are drawn i.i.d. $i = 1, \dots, n$. Consider the statistic

$$\delta(X) = \max\{X_1, \dots, X_n\}.$$

Its distribution is easily found to be

$$\mathbb{P}(\delta(X) \leq t) = \begin{cases} 0 & \text{if } t < 0 \\ \left(\frac{t}{\theta}\right)^n & \text{if } 0 \leq t \leq \theta \\ 1 & \text{if } t > \theta \end{cases}$$

with density $f(t) = \frac{n}{\theta^n} t^{n-1}$ for $t \in [0, \theta]$. We thus have

$$\mathbb{E}_\theta(\delta(X)) = \frac{n}{\theta^n} \int_0^\theta t^n dt = \frac{n}{n+1}\theta.$$

Similarly,

$$\mathbb{E}_\theta(\delta^2(X)) = \frac{n}{\theta^n} \int_0^\theta t^{n+1} dt = \frac{n}{n+2}\theta^2.$$

Clearly $\frac{n+1}{n}\delta(X)$ is an unbiased estimator. But instead consider all estimators of the form $a\delta(X)$, $a \in \mathbb{R}$. Find the value of a that gives the minimal risk with the quadratic loss function. We have

$$\begin{aligned} R(\theta, a\delta) &= \mathbb{E}_\theta(\theta - a\delta(X))^2 = \theta^2 - 2a\mathbb{E}_\theta(\delta(X)) + a^2\mathbb{E}_\theta\delta(X)^2 \\ &= \theta^2 - 2a\frac{n}{n+1}\theta^2 + a^2\frac{n}{n+2}\theta^2. \end{aligned}$$

This is a convex function optimum does not depend on θ and the minimum is easily found to be $a^* = \frac{n+2}{n+1}$. We have

$$R(\theta, \frac{n+1}{n}\delta) = \frac{\theta^2}{n(n+2)}, \quad R(\theta, \frac{n+2}{n+1}\delta) = \frac{\theta^2}{(n+1)^2}.$$

This shows that unbiased estimators are not always best if the goal is to minimize risk.

There are situations where certain decision rules can be disregarded. A decision rule $\delta(X)$ is **inadmissible** if there is some competing procedure δ' has uniformly lower risk, meaning:

- (a) $R(\theta, \delta') \leq R(\theta, \delta)$ for all $\theta \in \Theta$,
- (b) $R(\theta, \delta') < R(\theta, \delta)$ for at least one $\theta \in \Theta$.

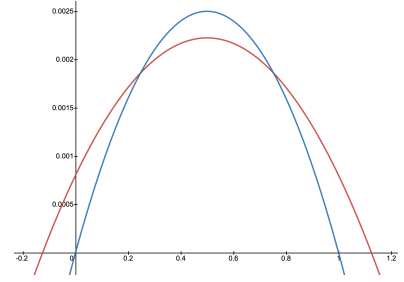


Figure 2.1: Risks functions for δ_1 (blue) and δ_3 (red) shown for $n = 100$.

Otherwise the decision rule is **admissible**. We treat admissibility in more detail in Section 2.4. In practice however, we often need to compare two admissible procedures and then there is no one obvious way to say that one is better than the other; more on that in Section 2.3.

Example 2.1.2 shows that a natural unbiased estimator may not be admissible. A slightly more artificial example follows.

Example 2.1.3. Continuing Example 2.1.1, consider $\delta_4(X) = X_1$, a rather silly estimator that uses only the first observation. Then $R(\theta, \delta_4) = \mathbb{E}_\theta(\theta - X_1)^2 = \theta(1 - \theta)$ which is always greater than $R(\theta, \delta_1)$ for all $\theta \in (0, 1)$ (unless $n = 1$ of course). Therefore δ_4 is inadmissible.

Given a sample $X = (X_1, \dots, X_n)$, define the **empirical risk** as

$$\widehat{R}_n(\theta, \delta) = \frac{1}{n} \sum_{i=1}^n L(\theta, \delta(X_i)). \quad (2.1)$$

A popular way of constructing an estimator is by minimizing the empirical risk, which is directly related to M-estimation discussed later in Section 8.3.

← Exercise 2.7.2

Remark 2.1.4. In the supervised learning set-up we have the training data $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$. Where $X_i \in \mathcal{X}$ are called features and $Y_i \in \mathcal{Y} \subset \mathbb{R}$ are called labels. Here the loss function $L(y, \delta(x))$ is defined on $\mathcal{Y} \times \mathcal{X}$. The (expected) risk is then defined through a double integral over $\mathcal{X} \times \mathcal{Y}$. In this context the empirical risk becomes

$$\widehat{R}_n(\delta) = \frac{1}{n} \sum_{i=1}^n L(Y_i, \delta(X_i)).$$

2.2 Randomized decision rules*

It can often be useful to consider randomized decision rules.

Definition 2.2.1. A **randomized decision rule** δ is a measurable mapping from \mathcal{X} to probability measures on \mathcal{A} , $x \mapsto \delta_x$. By this we mean that the function $x \mapsto \delta_x(B) \in [0, 1]$ is measurable for any fixed Borel subset $B \subseteq \mathcal{A}$. (such objects are also called probability kernels)

The idea is that for a fixed $X = x$ a random action A will be drawn from δ_x , $A|X = x \sim \delta_x$, that is, $\mathbb{P}(A \in B|X = x) = \delta_x(B)$ for every Borel set $B \subseteq \mathcal{A}$. The definition makes sure that the joint distribution probabilities of the the following form are well defined:

$$\mathbb{P}(A \in B, X \in U) = \int_U \int_B \delta_x(da) \mathbb{P}(dx) = \int_U \delta_x(B) \mathbb{P}(dx),$$

which they are because $\delta_x(B)$ is measurable in x .

Note that the marginal distribution of X and the conditional distribution of A given X naturally specify the joint distribution of (A, X) and so also the marginal distribution of A . For a randomized decision rule, we define the risk function as

$$R(\theta, \delta) = \mathbb{E}_\theta L(\theta, A) = \mathbb{E}_\theta (\mathbb{E}[L(\theta, A)|X]) = \iint L(\theta, a) \delta_x(da) P_\theta(dx). \quad (2.2)$$

Here the conditional expectation $\mathbb{E}(\cdot|X = x)$ is computed with respect to the conditional distribution δ_x of A given $X = x$. Denote by \mathcal{D} the set of all randomized decision rules such that $R(\theta, \delta) < +\infty$ for all $\theta \in \Theta$. The set \mathcal{D} contains all non-randomized decision rules \mathcal{D}_0 in which case to each x we assign the point mass at this point.

Example 2.2.2 (Statistical testing). Consider the problem of testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$. Here $\mathcal{A} = \{0, 1\}$. We consider here a randomized testing procedure where based on data $x \in \mathcal{X}$ we compute $\varphi(x) \in [0, 1]$ and let δ_x be the Bernoulli distribution with the success probability $\varphi(x)$. In testing problems we typically use the 0/1-loss:

$$L(\theta, a) = \mathbb{1}(a = 1, \theta \in \Theta_0) + \mathbb{1}(a = 0, \theta \in \Theta_1).$$

The **power function** assigns to θ the probability of rejecting H_0 :

$$\beta(\theta) := \mathbb{P}_\theta(A = 1) = \mathbb{E}[\mathbb{P}_\theta(A = 1|X)] = \mathbb{E}\varphi(X). \quad (2.3)$$

The risk function for the randomized rule defined above is

$$\begin{aligned} R(\theta, \delta) &= \mathbb{E}_\theta \mathbb{E}[L(\theta, A)|X] = \mathbb{E}_\theta \mathbb{P}(A = 1, \theta \in \Theta_0|X) + \mathbb{E}_\theta \mathbb{P}(A = 0, \theta \in \Theta_1|X) \\ &= \begin{cases} \mathbb{E}_\theta \varphi(X) & \text{if } \theta \in \Theta_0 \\ \mathbb{E}_\theta (1 - \varphi(X)) & \text{if } \theta \in \Theta_1 \end{cases} \\ &= \begin{cases} \beta(\theta) & \text{if } \theta \in \Theta_0 \\ 1 - \beta(\theta) & \text{if } \theta \in \Theta_1 \end{cases}. \end{aligned}$$

In standard statistical theory, randomized decision rules appear mostly in the context of statistical testing. It is important however to discuss briefly general advantages of considering them. For two randomized decision rules δ, δ' and $\lambda \in (0, 1)$ the convex combination $(1 - \lambda)\delta + \lambda\delta'$ is the randomized decision rule defined by

$$(1 - \lambda)\delta + \lambda\delta' = \begin{cases} \delta & \text{with probability } 1 - \lambda, \\ \delta' & \text{with probability } \lambda. \end{cases} \quad (2.4)$$

With this definition the set of randomized decision rules \mathcal{D} forms a convex set. The following proposition suggests that convexity will play an important role in our analysis; see also Remark 2.3.1 below.

Proposition 2.2.3. *The set \mathcal{D} of all randomized decision rules is convex. The function $\delta \mapsto R(\theta, \delta)$ is a linear function of $\delta \in \mathcal{D}$ for any fixed θ , that is,*

$$R(\theta, (1 - \lambda)\delta + \lambda\delta') = (1 - \lambda)R(\theta, \delta) + \lambda R(\theta, \delta')$$

for every $\lambda \in (0, 1)$ and $\delta, \delta' \in \mathcal{D}$.

Proof. The first statement is clear. For the second statement let A_0, A_1 be the random variables representing the randomized decisions δ and δ' . Their mixture is A and so A is equal to A_0 with probability $(1 - \lambda)$ and it is equal to A_1 with probability λ . We thus have

$$R(\theta, (1 - \lambda)\delta + \lambda\delta') = \mathbb{E}_\theta L(\theta, A) = (1 - \lambda)\mathbb{E}_\theta L(\theta, A_0) + \lambda\mathbb{E}_\theta L(\theta, A_1),$$

which is equal to $(1 - \lambda)R(\theta, \delta) + \lambda R(\theta, \delta')$ as claimed. \square

It may be surprising that linearity in Proposition 2.2.3 holds irrespective of the loss function. The following remark may help.

Remark 2.2.4. *It is important to note a subtlety in the above discussion. If δ, δ' are non-randomized decision rules then $(1 - \lambda)\delta + \lambda\delta'$ could mean two different objects. It could be either a randomized decision rule as defined in (2.4), but it could be also a non-randomized decision rule*

$$((1 - \lambda)\delta + \lambda\delta')(x) = (1 - \lambda)\delta(x) + \lambda\delta'(x).$$

For the last definition, there is no equivalent of Proposition 2.2.3 although convexity of \mathcal{D}_0 still holds as long as \mathcal{A} is convex. In the hypothesis testing case, $\mathcal{A} = \{0, 1\}$ is not convex.

2.3 Bayesian and minimax rules*

In the ideal situation we would be able to choose the best decision rule δ as the one that leads to the lowest risk. However the risk $R(\theta, \delta)$ of δ depends on θ and, in general, there is no total ordering on the risk functions for different θ ; c.f Figure 2.1. The only realistic solution is to define a functional on the space of risk functions and evaluate decision rules according to the value of this functional.

2.3.1 Definitions and basic properties

Two natural choices for functionals on the risk functions are the maximum risk and the Bayes risk. The maximum risk of δ is defined as the number

$$\bar{R}(\delta) = \sup_{\theta \in \Theta} R(\theta, \delta).$$

For a given prior distribution $\pi(\theta)$ on Θ , the corresponding **Bayes risk** is

$$r(\pi, \delta) = \int R(\theta, \delta)\pi(\theta)d\theta.$$

We say that a decision rule δ^* is a (randomized) **minimax rule** if

$$\bar{R}(\delta^*) = \inf_{\delta \in \mathcal{D}} \bar{R}(\delta).$$

A (randomized) **Bayes rule** with respect to prior π is any decision rule δ^* that satisfies

$$r(\pi, \delta^*) = \inf_{\delta \in \mathcal{D}} r(\pi, \delta).$$

When the infima are over \mathcal{D}_0 we call the corresponding decision procedures **non-randomized Bayesian** and **non-randomized minimax** rules respectively.

In practice we prefer non-randomized decision rules. The reason for considering randomized versions can be explained as follows.

Remark 2.3.1. By Proposition 2.2.3, $R(\theta, \delta)$ is linear (and so convex) in δ for any fixed θ . By Proposition B.2.5, $\bar{R}(\delta)$ is convex and $r(\pi, \delta)$ is linear in δ . It follows that both minimizing $\bar{R}(\delta)$ and $r(\pi, \delta)$ over \mathcal{D} is a (infinite dimensional) convex optimization problem. For non-randomized rules such favorable properties will hold only under special assumptions on the loss function.

The following result shows that there is no difference between randomized and non-randomized Bayes rules.

Theorem 2.3.2. Every Bayes rule δ^* satisfies $r(\pi, \delta^*) = \inf_{\delta \in \mathcal{D}_0} r(\pi, \delta)$.

Proof. We omit the formal proof. See Section 1.8 in ². □

2.3.2 Simple geometric insights

It is good to briefly discuss the geometric picture that drives our intuition. Consider the set of all finite mixtures of all non-randomized decision rules. It is natural to denote this set as $\text{conv}(\mathcal{D}_0)$. Although in general $\mathcal{D} \neq \text{conv}(\mathcal{D}_0)$ from the point of view of risk analysis often we can use these two sets interchangeably obtaining the same results; see Section 1.6 in ³ for a more careful discussion and further references. This insight makes Theorem 2.3.2 very natural to conjecture because $r(\pi, \delta)$ is a linear function of δ and we optimize it (essentially) over $\text{conv}(\mathcal{D}_0)$.

We next discuss a trivial instance of this when things can be explicitly depicted in two dimensions. Take $\Theta = \{0, 1\}$ the risk function $R(\theta, \delta)$ can be represented by a point $(R(0, \delta), R(1, \delta)) \in \mathbb{R}^2$. Define the risk set

$$\mathcal{R} = \{(R(0, \delta), R(1, \delta)) : \delta \in \text{conv}(\mathcal{D}_0)\}.$$

By Proposition 2.2.3, \mathcal{R} is a convex set, namely, if $y, y' \in \mathcal{R}$ (with the underlying δ, δ') then $(1 - \lambda)y + \lambda y'$ lies in \mathcal{R} as it is realized by the randomized rule $(1 - \lambda)\delta + \lambda\delta'$.

² Thomas S. Ferguson. *Mathematical statistics: A decision theoretic approach*. Probability and Mathematical Statistics, Vol. 1. Academic Press, New York-London, 1967

³ Thomas S. Ferguson. *Mathematical statistics: A decision theoretic approach*. Probability and Mathematical Statistics, Vol. 1. Academic Press, New York-London, 1967

Suppose that the considered set of non-randomized decision rules \mathcal{D}_0 is finite. In this case \mathcal{R} forms a polytope and the geometry of the Bayes rule is quite straightforward. A prior distribution π is simply a point $(\pi(0), \pi(1))$ with nonnegative coordinates that sum to 1.

Moreover,

$$r(\pi, \delta) = \pi(0)R(0, \delta) + \pi(1)R(1, \delta).$$

Thus a Bayes rule will be obtained as one of the minimizers of a given linear function on \mathcal{R} . The maximizers may be unique but it is not always the case. In this case, Theorem 2.3.2 holds trivially.

For the minimax rule note that

$$\bar{R}(\delta) = \max\{R(0, \delta), R(1, \delta)\}.$$

Draw the line $R(0, \delta) = R(1, \delta)$ and consider how the minimax risk can be improved above or below the line.

2.3.3 Finding the Bayes rule

We next explain how to find a Bayes rule. By Theorem 2.3.2, it is enough to minimize $r(\pi, \delta)$ over $\delta \in \mathcal{D}_0$. For a prior distribution $\pi(\theta)$ let $\pi(\theta|x)$ denote the posterior distribution. We write $p(x|\theta)$ be the density of the distribution \mathbb{P}_θ . By the Bayes theorem

$$\pi(\theta|x) = \frac{p(x|\theta)\pi(\theta)}{\int p(x|\theta)\pi(\theta)d\theta} = \frac{p(x|\theta)\pi(\theta)}{m(x)}, \quad (2.5)$$

where $m(x) = \int p(x|\theta)\pi(\theta)d\theta$. The **posterior risk** is defined as

$$r(\delta|x) := \int L(\theta, \delta(x))\pi(\theta|x)d\theta. \quad (2.6)$$

One of the main result on Bayesian rules is the following proposition.

Proposition 2.3.3. *The Bayes risk $r(\pi, \delta)$ for $\delta \in \mathcal{D}_0$ satisfies*

$$r(\pi, \delta) = \int r(\delta|x)m(x)dx.$$

Define δ^* pointwise, for every $x \in \mathcal{X}$, as

$$\delta^*(x) = \delta_0(x), \quad \text{where } r(\delta_0|x) = \inf_{\delta \in \mathcal{D}_0} r(\delta|x). \quad (2.7)$$

If $\delta^* \in \mathcal{D}_0$ (i.e. if δ^* is measurable) then, by construction, $r(\delta^*|x) = \inf_{\delta \in \mathcal{D}_0} r(\delta|x)$ for every $x \in \mathcal{X}$ and δ^* is a Bayes rule.

Proof. Using the Fubini's theorem we get

$$\begin{aligned} r(\pi, \delta) &= \int R(\theta, \delta)\pi(\theta)d\theta = \iint L(\theta, \delta(x))p(x|\theta)\pi(\theta)dx d\theta \\ &\stackrel{(2.5)}{=} \iint L(\theta, \delta(x))\pi(\theta|x)m(x)dx d\theta \\ &= \int \left(\int L(\theta, \delta(x))\pi(\theta|x)d\theta \right) m(x)dx \\ &= \int r(\delta|x)m(x)dx. \end{aligned}$$

It is also clear that δ^* minimizes the above integral over \mathcal{D}_0 :

$$\inf_{\delta \in \mathcal{D}_0} r(\pi, \delta) = \inf_{\delta \in \mathcal{D}_0} \int r(\delta|x)m(x)dx \geq \int (\inf_{\delta \in \mathcal{D}_0} r(\delta|x))m(x)dx = r(\pi, \delta^*).$$

If $\delta^* \in \mathcal{D}_0$ (c.f. Remark 2.3.4), the inequality becomes equality. By Theorem 2.3.2, minimizing over \mathcal{D}_0 is equivalent to minimizing over \mathcal{D} showing that δ^* is a Bayes rule. \square

Remark 2.3.4. The function δ^* defined in Proposition 2.3.3 will be however measurable in all practical situations. For more details see 4.

⁴ Lawrence D Brown and Roger Purves. Measurable selections of extrema. *The Annals of Statistics*, pages 902–912, 1973

In some situations the Bayes rule can be found explicitly.

Proposition 2.3.5. If $\Theta, \mathcal{A} \subseteq \mathbb{R}^d$ and $L(\theta, a) = \|\theta - a\|^2$ then the Bayes rule is

$$\delta^*(x) = \int \theta \pi(\theta|x)d\theta = \mathbb{E}(\theta|X = x).$$

Proof. By Proposition 2.3.3, the decision rule $\delta^* = \inf_{\delta \in \mathcal{D}_0} r(\delta|x)$ is a Bayes rule. In our case

$$r(\delta|x) = \int \|\theta - \delta(x)\|^2 \pi(\theta|x)d\theta.$$

For a fixed x , write $a = \delta(x)$ and minimize the above expression with respect to a . We have

$$\nabla_a \int \|\theta - a\|^2 \pi(\theta|x)d\theta = -2 \int (\theta - a) \pi(\theta|x)d\theta,$$

which is equal to the zero vector if and only if $a = \int \theta \pi(\theta|x)d\theta = \mathbb{E}(\theta|X = x)$. \square

Example 2.3.6. Given a sample X_1, \dots, X_n from $N(\theta, 1)$ consider the prior $\pi \sim N(0, \tau^2)$. The posterior is

$$N\left(\frac{n\tau^2}{n\tau^2 + 1} \bar{x}_n, \frac{\tau^2}{n\tau^2 + 1}\right).$$

Thus, for this prior and with the quadratic loss, the Bayes rule is

$$\delta^*(x_1, \dots, x_n) = \frac{n\tau^2}{n\tau^2 + 1} \bar{x}_n.$$

The risk function is

$$R(\theta, \delta^*) = \mathbb{E}_\theta(\theta - \frac{n\tau^2}{n\tau^2 + 1} \bar{X}_n)^2 = \left(\frac{1}{n\tau^2 + 1}\right)^2 \theta^2 + \frac{1}{n} \left(\frac{n\tau^2}{n\tau^2 + 1}\right)^2.$$

Since $\mathbb{E}\theta^2 = \tau^2$, the Bayes risk is

$$r(\pi, \delta^*) = \mathbb{E}_\theta(\theta - \frac{n\tau^2}{n\tau^2 + 1} \bar{X}_n)^2 = \left(\frac{1}{n\tau^2 + 1}\right)^2 \tau^2 + \frac{1}{n} \left(\frac{n\tau^2}{n\tau^2 + 1}\right)^2.$$

For comparison, consider $\delta_0(x_1, \dots, x_n) = \bar{x}_n$. We have $R(\theta, \delta_0) = \frac{1}{n}$ and so $r(\pi, \delta_0) = \frac{1}{n}$.

← Exercise 2.7.3

← Exercise 2.7.4

← Exercise 2.7.5

2.3.4 The link between minimax and Bayesian rules

The next result gives a connection between minimax and Bayesian rules.

Theorem 2.3.7. *Let δ^* be a Bayes rule for some prior π . Suppose that*

$$R(\theta, \delta^*) \leq r(\pi, \delta^*) \quad \text{for all } \theta \in \Theta. \quad (2.8)$$

Then δ^ is minimax.*

Here probably the most interesting case when this happens is when $R(\theta, \delta^*)$ is constant in θ . Compare this with the geometric picture presented in Section 2.3.2.

Proof. We prove the contrapositive statement. If δ^* is not minimax then there exists δ_0 such that $\bar{R}(\delta_0) < \bar{R}(\delta^*)$. However, the fact that $r(\pi, \delta_0) \leq \sup_{\theta} R(\theta, \delta_0)$ implies that

$$r(\pi, \delta_0) \leq \bar{R}(\delta_0) < \bar{R}(\delta^*) \stackrel{(2.8)}{\leq} r(\pi, \delta^*).$$

But this contradicts the fact that δ^* is a Bayesian rule. \square

To illustrate the utility of this theorem in actually finding minimax estimators consider the following example.

Example 2.3.8 (Bernoulli distribution). *Suppose $X \sim \text{Bern}(\theta)$, $\theta \in \Theta = [0, 1]$. Given any prior π over $[0, 1]$, we define $m_1 = \mathbb{E}(\theta)$ and $m_2 = \mathbb{E}(\theta^2)$. The frequentist risk for the squared loss is*

$$R(\theta, \delta) = \theta^2(1 + 2(\delta_0 - \delta_1)) + \theta(\delta_1^2 - \delta_0^2 - 2\delta_0) + \delta_0^2,$$

whereas the Bayes risk is

$$r(\pi, \delta) = m_2(1 + 2(\delta_0 - \delta_1)) + m_1(\delta_1^2 - \delta_0^2 - 2\delta_0) + \delta_0^2.$$

Note that it depends on π only through the first two moments! The Bayes decision rule is found minimizing with respect to δ_0, δ_1 :

$$\delta_0 = \frac{m_1 - m_2}{1 - m_1}, \quad \delta_1 = \frac{m_2}{m_1}. \quad (2.9)$$

To satisfy the constant risk property, we note that $R(\theta, \delta)$ does not depend on θ if and only if $\delta_1 - \delta_0 = \frac{1}{2}$ and $\delta_1^2 - \delta_0^2 - 2\delta_0 = 0$, or equivalently $\delta_0 = \frac{1}{4}$, $\delta_1 = \frac{3}{4}$. Solving for m_1, m_2 in (2.9) yields the solution

$$m_1^* = \frac{1}{2}, \quad m_2^* = \frac{3}{8}.$$

This, if π is such that $\mathbb{E}\theta = \frac{1}{2}$ and $\mathbb{E}\theta^2 = \frac{3}{8}$ and $\delta_0^ = \frac{1}{4}$, $\delta_1^* = \frac{3}{4}$, then δ^* is a Bayesian rule with respect to π and inequality (2.8) holds. By Theorem 2.3.7, δ^* is also minimax.*

Note that if $\pi = \text{Beta}(1/2, 1/2)$ then θ has these moments. The associated Bayes (and hence minimax) risk is $\frac{1}{16}$. These calculations can be generalized to the case of n i.i.d. samples X_1, \dots, X_n from the Bernoulli distribution. We leave it as an exercise.

Consider the following slight generalization of the minimax rule: allow nature to choose a distribution π over Θ . In this case our goal is to minimize

$$\bar{r}(\delta) := \sup_{\pi} r(\pi, \delta).$$

In this context it is useful to consider the notion of a least favourable prior. A prior π^* is called **least favorable** if

$$\inf_{\delta \in \mathcal{D}} r(\pi^*, \delta) = \sup_{\pi} \inf_{\delta \in \mathcal{D}} r(\pi, \delta).$$

← Exercise 2.7.6

2.3.5 Minimax theorem*

The minimax decision problem has a natural game-theoretic interpretation. In particular, a two-player game is played between nature and the statistician, with nature picking the prior π , and the statistician choosing a (possibly randomized) decision rule δ . Then the statistician pays to nature the amount $r(\pi, \delta)$. Note that nature gains the same amount the statistician loses, so that the game is zero-sum. The following two quantities are important

$$L^* = \sup_{\pi} \inf_{\delta} r(\pi, \delta)$$

$$U^* = \inf_{\delta} \sup_{\pi} r(\pi, \delta)$$

called respectively the lower and the upper values of the game. These quantities have the following interpretation: U^* is the amount the player pays when he is told what distribution nature chooses before he chooses δ . Conversely, L^* is the amount the player pays when nature is told the player's strategy δ before it chooses π .

← Exercise 2.7.7

Theorem 2.3.9 (von Neumann). *Suppose that*

1. *The parameter space $\Theta = \{\theta_1, \dots, \theta_k\}$ is finite, and*
2. *The risk set*

$$\mathcal{R} = \{y \in \mathbb{R}^k : y_i = R(\theta_i, \delta) \text{ for some } \delta \in \mathcal{D}\}$$

is closed and lies in the nonnegative orthant.

Then

- (i) *The game has the value $L^* = U^*$.*

(ii) There exists a probability vector $\pi \in \mathbb{R}^k$ that is a least favourable prior.

Proof. In Exercise 2.7.7 we show that $L^* \leq U^*$ so in order to establish

(i) we need to prove that $L^* \geq U^*$. Given $\alpha \in \mathbb{R}$, define the lower-rectangular set of the form

$$B_\alpha = \{y \in \mathbb{R}^k : y_i \leq \alpha, \forall i = 1, \dots, k\}.$$

In addition, define

$$\gamma = \arg \inf_{\alpha \in \mathbb{R}} \{B_\alpha \cap \mathcal{R} \neq \emptyset\}.$$

By definition of γ , for each $n \in \mathbb{N}$ there exists a (randomized) decision rule δ_n , such that

$$R(\theta_j, \delta_n) \leq \gamma + \frac{1}{n}, \quad \text{for all } j = 1, \dots, k.$$

Therefore, for any prior π on $\Theta = \{\theta_1, \dots, \theta_k\}$, we have

$$r(\pi, \delta_n) \leq \gamma + \frac{1}{n}.$$

Taking a supremum over the choice of priors yields that $\sup_\pi r(\pi, \delta_n) \leq \gamma + \frac{1}{n}$ and hence for $n \in \mathbb{N}$

$$\inf_\delta \sup_\pi r(\pi, \delta) \leq \gamma + \frac{1}{n}.$$

This inequality holds for every $n \in \mathbb{N}$ and taking the limit gives $U^* \leq \gamma$.

Now we are going to use the Separating Hyperplane Theorem B.1.3 to construct a vector $\pi = (\pi_1, \dots, \pi_k)$ that can be viewed as a least favourable prior, therefore establishing (ii). It will also show that $L^* \geq \gamma$ establishing part (i).

Consider again the lower rectangle B_γ . Observe that its interior $\text{int}(B_\gamma)$ and the risk set \mathcal{R} are two disjoint convex sets in \mathbb{R}^k . Consequently, the separating hyperplane theorem guarantees existence of some non-zero vector $\pi \in \mathbb{R}^k$ and constant c such that

$$\langle \pi, y \rangle \geq c \quad \text{for all } y \in \mathcal{R},$$

$$\langle \pi, y \rangle \leq c \quad \text{for all } y \in \text{int}(B_\gamma).$$

We claim that $\pi \geq 0$. This can be proven by contradiction. Suppose $\pi_i < 0$ for some i . Directly by definition of the set B_γ , we can construct a sequence of y vectors such that $y_i \rightarrow -\infty$ and $y_j = 0$ for $j \neq i$, while still staying in $\text{int}(B_\gamma)$. However, this yields a sequence such that $\langle \pi, y \rangle$ becomes indefinitely large contradicting the separation statement (these sequence ends up on the wrong side of the hyperplane $\langle \pi, y \rangle = c$). Thus we must have $\pi \geq 0$. Since $\pi \neq 0$ we can normalize it to sum to one, so it can be interpreted as a valid prior.

Consider the vector $x^* := \gamma \mathbf{1}$. Since x^* lies in the closure of $\text{int}(B_\gamma)$, we must have

$$\langle \pi, x^* \rangle = \gamma \leq c,$$

where we have used normalization property of π . Now letting δ be an arbitrary decision rule with risk vector $z \in \mathbb{R}^k$ such that $z_i = R(\theta_i, \delta)$ for all i we have

$$r(\pi, \delta) = \langle \pi, z \rangle \geq c \geq \gamma.$$

Since δ was arbitrary, $\inf_\delta r(\pi, \delta) \geq \gamma$ and, in consequence, $L^* \geq \gamma$, which completes the proof of part (i). Furthermore, the vector π that we constructed is the least favourable prior of part (ii). \square

2.4 Admissibility and Rao-Blackwell

It is clear that we can discard inadmissible rules from our analysis. Thus it is useful to know that the procedure we analyse is admissible. It should be clear that minimax procedures are admissible.

Theorem 2.4.1. *If a Bayes rule for π is essentially unique (in the sense of measurable functions) then δ is admissible.*

Proof. Suppose there exists δ_0 such that $R(\theta, \delta_0) \leq R(\theta, \delta)$ for all $\theta \in \Theta$. Then $r(\pi, \delta_0) \leq r(\pi, \delta)$ and so, if δ is Bayes, then δ_0 is also Bayes. By our assumption, $\delta = \delta_0$ almost surely. In consequence, for every θ

$$\mathbb{E}_\theta L(\theta, \delta(X)) = \mathbb{E}_\theta L(\theta, \delta_0(X))$$

and so the risk functions are equal implying that δ_0 does not strictly dominate δ . \square

In case, we cannot assure uniqueness, it is still useful to provide some sufficient conditions on admissibility.

We now revisit the Rao-Blackwell theorem in the language of statistical decision theory. It allows for a simple general construction of decision rules that dominate a given decision rule δ . Recall that a statistic $T = T(X)$ is sufficient for θ if the conditional distribution of X given T does not depend on θ . Given any rule, we can define

$$\eta(T) = \mathbb{E}[\delta(X)|T].$$

Since T is sufficient, $\mathbb{E}[\delta(X)|T]$ does not depend on θ and so, $\eta(T)$ is a valid statistic. The following classical result explains why $\eta(T)$ may be preferred over $\delta(X)$.

Theorem 2.4.2 (Rao-Blackwell). *Let T be a sufficient statistic for $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$, let δ be a decision rule, and define $\eta(T) = \mathbb{E}[\delta(X)|T]$. If $\theta \in \Theta$, $R(\theta, \delta) < +\infty$, and $L(\theta, a)$ is convex in a , then*

$$R(\theta, \eta) \leq R(\theta, \delta).$$

← Exercise 2.7.8

← Exercise 2.7.9

← Exercise 2.7.10

Proof. Using the Jensen's inequality (c.f Theorem B.2.3) for the convex function $f(a) = L(\theta, a)$ and with the conditional expectation $\mathbb{E}(\cdot|T = t)$ we get

$$L(\theta, \eta(T)) = L(\theta, \mathbb{E}(\delta(X)|T)) \leq \mathbb{E}(L(\theta, \delta(X))|T). \quad (2.10)$$

Taking expectations on both sides we conclude $R(\theta, \eta(T)) \leq R(\theta, \delta(X))$. \square

This theorem states that, as long as the loss function is convex, every estimator that depends on the data not through the sufficient statistics can be improved. In consequence, in our study of optimal procedures, we can always focus on procedures that are based on the sufficient statistics.

Remark 2.4.3. In Theorem 2.4.2 we can also show that if $L(\theta, a)$ is strictly convex in a , the inequality will be strict unless $\delta(X) = \eta(T)$ a.e.. A rough argument goes as follows: If f is strictly convex then, by Theorem B.2.3, the inequality in (2.10) becomes strict unless $\delta(X)$ becomes constant a.s. after conditioning on the event $\{T = t\}$ (for every t). This implies that $\delta(X)$ must be a function of T and so $\delta(X) = \mathbb{E}(\delta(X)|T)$ a.s. If this inequality is strict then it remains strict after taking the expectation.

← Exercise 2.7.11

Example 2.4.4. Suppose $X_i \sim U(0, \theta)$, $i = 1, \dots, n$. Consider $X_{(n)} := \max\{X_1, \dots, X_n\}$. The fact that $X_{(n)}$ is a sufficient statistics follows from the Fisher-Neyman factorization theorem. Indeed, denote $X_{(1)} := \min\{X_1, \dots, X_n\}$, then the density of the sample $\mathbf{x} = (x_1, \dots, x_n)$ is

$$f_{\theta}(\mathbf{x}) = \frac{1}{\theta^n} \mathbb{1}\{x_{(1)} \geq 0\} \mathbb{1}\{x_{(n)} \leq \theta\},$$

which is equal to $h(\mathbf{x})g_{\theta}(x_{(n)})$ with $h(\mathbf{x}) = \mathbb{1}\{x_{(1)} \geq 0\}$ and $g_{\theta}(x_{(n)}) = \frac{1}{\theta^n} \mathbb{1}\{x_{(n)} \leq \theta\}$.

Consider now an unbiased estimator of θ given as

$$\delta(X) = \frac{2}{n} \sum_{i=1}^n X_i.$$

Under the quadratic loss

$$R(\theta, \delta) = \text{var}_{\theta}(\delta(X)) = \frac{4}{n} \text{var}_{\theta}(X_1) = \frac{4}{n} \frac{\theta^2}{12} = \frac{\theta^2}{3n}.$$

Consider now the rao-blackwellized version of δ defined as $\eta(X_{(n)}) = \mathbb{E}[\delta(X)|X_{(n)}]$. By Theorem 2.4.2 it dominates δ . To compute this new estimator explicitly, fix i and note that, for every $t > 0$

$$\mathbb{E}[X_i|X_{(n)} = t, X_i < t] = \mathbb{E}[X_i|X_i < t] = \frac{t}{2}.$$

The first equation follows because $X_{(n)}$ is independent of X_i conditionally on the event $\{X_i < X_{(n)}\}$ (in this case the maximum is a function of the remaining variables). The second equality follows simply because the distribution of X_i conditionally on $X_i < t$ is $U(0, t)$, which follows by standard calculations. We have shown that

$$\mathbb{E}[X_i | X_{(n)}, X_i < X_{(n)}] = \frac{X_{(n)}}{2}.$$

Trivially,

$$\mathbb{E}[X_i | X_{(n)}, X_i = X_{(n)}] = X_{(n)}.$$

The probability of the event $\{X_i = X_{(n)}\}$ is $\frac{1}{n}$ and thus

$$\mathbb{E}[X_i | X_{(n)}] = \frac{1}{n} \mathbb{E}[X_i | X_{(n)}, X_i = X_{(n)}] + \frac{n-1}{n} \mathbb{E}[X_i | X_{(n)}, X_i < X_{(n)}] = \frac{n+1}{2n} X_{(n)},$$

which gives

$$\eta(X_{(n)}) = \frac{2}{n} \sum_{i=1}^n \mathbb{E}[X_i | X_{(n)}] = \frac{2n}{n} \frac{n+1}{2n} X_{(n)} = \frac{n+1}{n} X_{(n)}.$$

In Example 2.1.2 we showed that this is an unbiased estimator of θ . We showed that its risk satisfies $R(\theta, \eta) = \frac{\theta^2}{n(n+2)}$. If $n \geq 2$ we then indeed have $R(\theta, \eta) \leq R(\theta, \delta)$ confirming the Rao-Blackwell theorem. Note however that Example 2.1.2 also showed that $\eta(X_{(n)})$ is not admissible!

← Exercise 2.7.11

Theorem 2.4.5. Let X_1, \dots, X_n be a random sample from $N(\theta, 1)$. The sample mean \bar{X}_n is an admissible estimator of θ under the squared loss.

We are not going to prove this result as it is quite technical. The idea follows from the fact that \bar{X}_n is a limit of Bayes rules for priors $\pi_k = N(0, k)$. Indeed, in this case, by Example 2.3.6, the posterior is $N(\frac{k^2}{k^2+1} \bar{x}_n, \frac{k}{k^2+1})$ and so it concentrates around \bar{x}_n for large k . The Bayes rule is $\frac{k^2}{k^2+1} \bar{x}_n$. Since each of these is admissible, with a bit of more care we can argue that the same holds in the limit. More generally, it can be shown that every admissible procedure is a limit of Bayes procedures; see Theorem 3.40⁵.

⁵ Mark J. Schervish. *Theory of statistics*. Springer Series in Statistics. Springer-Verlag, New York, 1995

2.5 Stein's paradox

Consider a sample $X^{(1)}, \dots, X^{(n)} \sim N_d(\mu, \Sigma)$ with Σ known. Here we assume $\Sigma = \sigma^2 I_d$. In this case $\epsilon := X - \mu \sim N_d(0, \sigma^2 I_d)$. Using the square loss $L(\mu, a) = \|\mu - a\|^2$ we easily show that the risk of any estimator $\hat{\mu}_n$ admits the following decomposition

$$R(\mu, \hat{\mu}_n) = \mathbb{E} \|\hat{\mu}_n - \mathbb{E}[\hat{\mu}_n]\|^2 + \|\mathbb{E}[\hat{\mu}_n] - \mu\|^2.$$

The first term is the variance of $\hat{\mu}_n$ and the second term is related to its bias. The MLE estimator of μ is the sample average \bar{X}_n . Its bias is

This is called the Gaussian sequence model and will reappear in this lecture.

zero and its variance is

$$\mathbb{E} \left\| \frac{1}{n} \sum_i (X^{(i)} - \mu) \right\|^2 = \frac{1}{n} \mathbb{E} \|\epsilon\|^2 = \sigma^2 \frac{d}{n}. \quad (2.11)$$

For simplicity take $\sigma^2 = 1$. In the special case when $n = 1$, the MLE is X which is an unbiased estimator with large variance d . Consider first as an alternative a simple linear estimator $\hat{\mu}_C = CX$ with $C = \text{diag}(c)$ diagonal and $c = (c_1, \dots, c_d)$. In this case

$$R(\mu, \hat{\mu}_C) = \sum_{i=1}^d (1 - c_i)^2 \mu_i^2 + \sum_{i=1}^d c_i^2.$$

Suppose we restrict ourselves to the hyperrectangular model class $|\mu_i| \leq \tau_i$. In this case we can easily find that the minimax risk is

$$\inf_c \sup_{-\tau \leq \mu \leq \tau} R(\mu, \hat{\mu}_C) = \inf_c \sum_{i=1}^d (1 - c_i)^2 \tau_i^2 + \sum_{i=1}^d c_i^2 = \sum_{i=1}^d \frac{\tau_i^2}{1 + \tau_i^2} < d.$$

This computation shows that for sparse model classes diagonal estimators $\hat{\mu}_C$ may strictly dominate the MLE if c is chosen carefully.

In this section we show a famous surprising result that the observation $X = (X_1, \dots, X_d)$ is not an admissible estimator for the parameter $\mu = (\mu_1, \dots, \mu_d)$ unless $d \leq 2$. We start introducing the Stein's Unbiased Risk Estimates (SURE).

2.5.1 Stein's Unbiased Risk Estimates (SURE)

For a function $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$ denote by $Jh(x)$ the Jacobian of h at $x \in \mathbb{R}^d$, that is, the matrix of partial derivatives with

$$(Jh(x))_{ij} = \frac{\partial h_i}{\partial x_j}.$$

The goal of this section is to prove the following result.

Proposition 2.5.1. [Stein's Unbiased Risk Estimates (SURE)] Let X_1, \dots, X_d be independent, $X_i \sim N(\mu_i, 1)$. Consider an estimator $\hat{\mu}(X)$ of $\mu = (\mu_1, \dots, \mu_d)$ and let

$$h(x) = x - \hat{\mu}(x). \quad (2.12)$$

Suppose that $h(x)$ satisfies

- (i) h is differentiable,
- (ii) $\mathbb{E}_\mu \|Jh(X)\| < \infty$.

Define

$$\hat{R} = d + \|h(X)\|^2 - 2\text{tr}(Jh(X)).$$

Then

$$R(\mu, \hat{\mu}) = \mathbb{E}_\mu \|\hat{\mu}(X) - \mu\|^2 = \mathbb{E}_\mu \hat{R}.$$

Remark 2.5.2. Proposition 2.5.1 is true if condition (i) is replaced with (i') h is weakly differentiable. A primary application being soft-thresholding when

$$h_i(x) = h_i(x_i) = x_i - \begin{cases} x_i - \lambda & x_i > \lambda \\ 0 & |x_i| \leq \lambda \\ x_i + \lambda & x_i < -\lambda. \end{cases}$$

Here $h_i(x)$ is weakly differentiable in the sense that whenever $-\infty < a \leq b < \infty$ there exists $h'_i(x_i)$ such that

$$\int_a^b h'_i(x_i) dx_i = h_i(b) - h_i(a).$$

To prove Proposition 2.5.1, we start with the Stein's lemma:

Lemma 2.5.3 (Stein's lemma). Let $X \sim N(\mu, \sigma^2)$ and let $h : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable with $\mathbb{E}|h'(X)| < \infty$. Then

$$\mathbb{E}[(X - \mu)h(X)] = \sigma^2 \mathbb{E}h'(X).$$

Proof. First assume $\mu = 0, \sigma^2 = 1$. In this case we equivalently show that $\mathbb{E}(Xh(X)) = \mathbb{E}h'(X)$. Without loss of generality assume $h(0) = 0^6$. The proof is an application of integration by parts. We have

$$(h(x)e^{-x^2/2})' = h'(x)e^{-x^2/2} - xh(x)e^{-x^2/2}$$

and so

$$0 = \mathbb{E}h'(X) - \mathbb{E}(Xh(X)),$$

which establishes the result. The fact that we get zero on the left follows from the fact that $\lim_{x \rightarrow \pm\infty} h(x)e^{-x^2/2} = 0$, which can be argued since $\mathbb{E}|h'(X)| < +\infty$. Indeed, first note that

$$|h'(t)|e^{-x^2/2}\mathbb{1}_{[0,x]}(t) \leq |h'(t)|e^{-t^2/2}$$

and the dominating function on the right is integrable by assumption. We also note that the smaller function goes to zero as $x \rightarrow \pm\infty$. Thus, we can use the dominating convergence theorem

$$\begin{aligned} \lim_{x \rightarrow \pm\infty} h(x)e^{-x^2/2} &= \lim_{x \rightarrow \pm\infty} \int_0^x h'(t)e^{-x^2/2} dt = \lim_{x \rightarrow \pm\infty} \int_{-\infty}^{\infty} h'(t)e^{-x^2/2}\mathbb{1}_{[0,x]}(t) dt \\ &= \int_{-\infty}^{\infty} \lim_{x \rightarrow \pm\infty} h'(t)e^{-x^2/2}\mathbb{1}_{[0,x]}(t) dt = 0. \end{aligned}$$

This establishes the result in the standard normal case. For general μ and σ define $Z = (X - \mu)/\sigma \sim N(0, 1)$. Define $\tilde{h}(z) = h(\mu + \sigma z)$. We have

$$\begin{aligned} \mathbb{E}[(X - \mu)h(X)] &= \sigma \mathbb{E}(Z h(\mu + \sigma Z)) = \sigma \mathbb{E}(Z \tilde{h}(Z)) \\ &= \sigma \mathbb{E}\tilde{h}'(Z) = \sigma^2 \mathbb{E}h'(\mu + \sigma Z) = \sigma^2 \mathbb{E}h'(X), \end{aligned}$$

where moving from the first to the second line we used the case $\mu = 0, \sigma = 1$ proved earlier. \square

⁶ Both sides will not change if we replace h with $h - h(0)$.

We need another technical lemma.

Lemma 2.5.4. *Let $X = (X_1, \dots, X_d)$ be a random vector with independent entries, $X_i \sim N(\mu_i, 1)$. If $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$ satisfies $\mathbb{E}\|Jh(X)\| < \infty$. Then*

$$\mathbb{E}((X - \mu)^\top h(X)) = \mathbb{E}\text{tr}(Jh(X))$$

Proof. Denote by $X_{\setminus i}$ the vector X with X_i removed. Using Lemma 2.5.3, for every $i = 1, \dots, d$

$$\begin{aligned} \mathbb{E}((X_i - \mu_i)h_i(X)) &= \mathbb{E}\left(\mathbb{E}[(X_i - \mu_i)h_i(X)|X_{\setminus i}]\right) \\ &= \mathbb{E}\left(\mathbb{E}\left[\frac{\partial h_i(X)}{\partial x_i} \middle| X_{\setminus i}\right]\right) = \mathbb{E}\left(\frac{\partial h_i(X)}{\partial x_i}\right). \end{aligned}$$

Summing over i we get the result. \square

Now Proposition 2.5.1 follows easily.

Proof of Proposition 2.5.1. We have

$$\begin{aligned} R(\mu, \hat{\mu}) &= \mathbb{E}((\hat{\mu}(X) - \mu)^\top (\hat{\mu}(X) - \mu)) \\ &= \mathbb{E}((\hat{\mu}(X) - X + X - \mu)^\top (\hat{\mu}(X) - X + X - \mu)) \\ &= \mathbb{E}\|h(X)\|^2 - 2\mathbb{E}((X - \mu)^\top h(X)) + \mathbb{E}\|X - \mu\|^2 \\ &\stackrel{\text{Lem 2.5.4}}{=} \mathbb{E}\|h(X)\|^2 - 2\mathbb{E}(\text{tr}(Jh(X))) + \mathbb{E}\|X - \mu\|^2 \\ &= \mathbb{E}\hat{R}, \end{aligned}$$

where the last equality follows because $\mathbb{E}\|X - \mu\|^2 = d$, which follows because $X - \mu$ is standard normal. \square

A natural estimator of μ is X and it has constant risk $R(\mu, X) = \mathbb{E}\|X - \mu\|^2 = d$. Although it is unbiased, the variance is large if d is large. The James-Stein estimator of $\mu = (\mu_1, \dots, \mu_d)$ is defined as

$$\delta_{\text{JS}}(X) = \left(1 - \frac{d-2}{\|X\|^2}\right)X.$$

Theorem 2.5.5. *The risk of the James-Stein estimator is*

$$R(\mu, \delta_{\text{JS}}) = d - \mathbb{E}_\mu \left(\frac{d-2}{\|X\|} \right)^2.$$

In particular, the natural estimator is not admissible if $d \geq 3$.

Proof. For the James-Stein estimator the function h in Proposition 2.5.1 is

$$h(x) = \frac{d-2}{\|x\|^2}x$$

and so

$$\hat{R} = d + \frac{(d-2)^2}{\|X\|^2} - 2\text{tr}(Jh(X)).$$

The diagonal entries of the Jacobian have a simple form

$$\frac{\partial h_i}{\partial x_i} = \frac{d-2}{\|x\|^2} - 2 \frac{d-2}{\|x\|^4} x_i^2$$

giving

$$\text{tr}(Jh(X)) = \left(\frac{d-2}{\|X\|} \right)^2 \quad \text{and} \quad \hat{R} = d - \left(\frac{d-2}{\|X\|} \right)^2.$$

By Proposition 2.5.1, if $d \geq 3$

$$R(\mu, \delta_{JS}) = \mathbb{E}_\mu \hat{R} = d - \mathbb{E}_\theta \left(\frac{d-2}{\|X\|} \right)^2 < d = R(\mu, X)$$

proving inadmissibility. \square

← Exercise 2.7.12

It is interesting to study the risk of a general linear estimator $\hat{\mu}_C = CX$ for an arbitrary square matrix C . To use SURE note that $h(X) = (I_d - C)X$ and so $Jh(x) = (I_d - C)$ giving $\text{tr}(Jh(x)) = d - \text{tr}(C)$. Therefore

$$R(\mu, \hat{\mu}_C) = \mathbb{E}_\mu [d + \|(I_d - C)X\|^2 - 2d + 2\text{tr}(C)] = \mathbb{E} \|(C - I_d)X\|^2 - d + 2\text{tr}(C). \quad (2.13)$$

We formulate the following result without a proof.

Proposition 2.5.6. *The linear estimator $\hat{\mu}_C = CX$ is admissible if and only if*

- (i) C is symmetric,
- (ii) the eigenvalues satisfy $0 \leq \rho_i(C) \leq 1$,
- (iii) $\rho_i(C) = 1$ for at most two i .

← Exercise 2.7.13

In many situations in the expression (2.13) $C = C(\lambda)$ depends on a regularization parameter λ so one could find optimal λ^* by minimizing the MSE over λ . One example is given by ridge regression. In this case for some $Z \in \mathbb{R}^{d \times p}$

$$C = Z(Z^\top Z + \lambda I_d)^{-1} Z^\top \quad (2.14)$$

and the corresponding estimator is obtained as $\hat{\mu}_Z = Z\hat{\beta}$ where $\hat{\beta} \in \mathbb{R}^p$ is obtained by solving

← Exercise 2.7.14

$$\min_{\beta} \frac{1}{2} \|X - Z\beta\|^2 + \frac{\lambda}{2} \|\beta\|^2, \quad \lambda > 0.$$

2.6 Minimizing risk under constraints

Earlier we saw that finding a decision rule δ that minimizes the risk $R(\theta, \delta)$ uniformly over θ is impossible. In Section 2.3 we saw two common strategies for introducing a global risk and finding optimal decision rules. An alternative approach is to introduce a reasonable constraint on the set of decision rules one is willing to consider. In this case, it can be possible to find a (constrained) δ for which $R(\theta, \delta)$ is minimized uniformly over θ . It is not so surprising that the interesting constraints are often convex. We review some of the most popular.

2.6.1 Unbiasedness constraints

An estimator is called unbiased if $\mathbb{E}_\theta \delta(X) = \theta$ for all θ . Sometimes the following definition is used instead.

Definition 2.6.1. For a loss $L(\theta, a)$ a decision rule δ is **unbiased with respect to L** if

$$\mathbb{E}_\theta(L(\theta', \delta(X))) \geq \mathbb{E}_\theta(L(\theta, \delta(X))) = R(\theta, \delta) \quad \text{for all } \theta, \theta' \in \Theta.$$

(Note that on both sides the expectation is taken with respect to \mathbb{P}_θ .)

Exercise 2.7.15 explores the connection between unbiased decision rules and unbiased estimation. Exercise 2.7.17 partially motivates the more general definition.

← Exercise 2.7.15

← Exercise 2.7.17

Proposition 2.6.2. The set of all (randomized) unbiased decision rules is convex.

Proof. In exactly the same way as in the proof of Theorem 2.2.3 we can show that $\mathbb{E}_\theta L(\theta', \delta)$ is linear in δ for any θ' . Let $\lambda \in (0, 1)$ and suppose that both δ and δ' are unbiased. It follows that

$$\begin{aligned} \mathbb{E}_\theta L(\theta', (1 - \lambda)\delta + \lambda\delta') &= (1 - \lambda)\mathbb{E}_\theta L(\theta', \delta) + \lambda\mathbb{E}_\theta L(\theta', \delta') \\ &\geq (1 - \lambda)\mathbb{E}_\theta L(\theta, \delta) + \lambda\mathbb{E}_\theta L(\theta, \delta') \\ &= \mathbb{E}_\theta L(\theta, (1 - \lambda)\delta + \lambda\delta') \end{aligned}$$

proving that $(1 - \lambda)\delta + \lambda\delta'$ is unbiased. \square

Though unbiasedness is more general, we shall focus here on the estimation problem. The next result says we need not look at Bayes estimators in this context, because (except in weird cases) they cannot be unbiased.

Proposition 2.6.3. No unbiased estimator $\delta(X)$ of $\theta \in \Theta \subseteq \mathbb{R}^d$ can be a Bayes estimator (under the square loss) unless the prior π satisfies $\pi(\{\theta : R(\theta, \delta) = 0\}) = 1$.

Proof. Suppose δ is a Bayes rule (under square loss with respect to π) and is unbiased. By Proposition 2.3.5 and the fact that δ is unbiased:

$$\delta(X) = \mathbb{E}(\theta|X) \quad \text{and} \quad \theta = \mathbb{E}(\delta(X)|\theta).$$

Then depending on the order in which we condition, we get

$$\mathbb{E}(\theta^\top \delta(X)) = \begin{cases} \mathbb{E}[\theta^\top \mathbb{E}(\delta(X)|\theta)] = \mathbb{E}(\theta^\top \theta) \\ \mathbb{E}[\delta(X)^\top \mathbb{E}(\theta|X)] = \mathbb{E}(\delta(X)^\top \delta(X)) \end{cases}.$$

Therefore $\mathbb{E}(\theta^\top \delta(X)) = \mathbb{E}(\theta^\top \theta) = \mathbb{E}(\delta(X)^\top \delta(X))$ and

$$\mathbb{E}\|\delta(X) - \theta\|^2 = \mathbb{E}(\delta(X)^\top \delta(X)) - 2\mathbb{E}(\theta^\top \delta(X)) + \mathbb{E}(\theta^\top \theta) = 0.$$

Since $r(\pi, \delta) = \mathbb{E}\|\delta(X) - \theta\|^2$ (with the expectation both with respect to X and θ) we get that $r(\pi, \delta) = 0$. But the Bayes risk also satisfies $r(\pi, \delta) = \int R(\theta, \delta) \pi(\theta) d\theta$. Since $R(\theta, \delta) \geq 0$ for all θ , the only way the π -integral can be zero is if π assigns probability 1 to the set of θ where $R(\theta, \delta)$ vanishes. This proves the claim. \square

Next we formulate a very powerful result, which states that there is a rule that uniformly minimizes risk over the unbiased estimators and, moreover, gives easily verifiable sufficient conditions to identify this best estimator. We first recall the notion of completeness.

Definition 2.6.4. A statistic T is called **complete** for the model $\mathcal{P} = \{\mathbb{P}_\theta\}$, if for every measurable function g , if $\mathbb{E}_\theta g(T) = 0$ for all θ then $g(T) = 0$ almost surely.

Although we did not specify this explicitly, a sufficient statistics in a regular exponential family is always complete.

Theorem 2.6.5 (Lehmann-Scheffe). Let $X \sim \mathbb{P}_\theta$ and suppose that T is a complete sufficient statistic. Suppose the goal is to estimate θ under convex loss, and that an unbiased estimator exists⁷. Then there exists an essentially unique unbiased estimator that is a function of T and uniformly minimizes the risk.

⁷ Note that in this result the estimator is unbiased in the classical sense but the loss function is general.

Proof. We first show that if δ is an unbiased estimator that uniformly minimizes risk then, without loss of generality, we can assume it is a function of T . Let δ be an unbiased estimator and define its rao-blackwellized version $\eta(T) = \mathbb{E}(\delta(X)|T)$ as in Theorem 2.4.2. Unbiasedness of δ gives

$$\theta = \mathbb{E}_\theta \delta(X) = \mathbb{E}_\theta [\mathbb{E}(\delta|T)] = \mathbb{E}_\theta [\eta(T)]$$

and so $\eta(T)$ is unbiased too. Moreover, by Theorem 2.4.2, $R(\theta, \eta) \leq R(\theta, \delta)$ for all θ .

Now we will show that the estimator $\eta(T)$ is essentially unique. Suppose $\eta^*(T)$ is also unbiased. Then

$$\mathbb{E}_\theta(\eta(T) - \eta^*(T)) = 0 \quad \text{for all } \theta \in \Theta$$

and by completeness $\eta(T) = \eta^*(T)$ a.s. showing that $\eta(T)$ is essentially unique. \square

Note that the proof also shows how such an estimator can be obtained. We simply start with an unbiased estimator and Rao-Blackwellize it. In Example 2.4.4 we essentially followed this construction. The estimator $\delta(X) = \frac{2}{n} \sum_{i=1}^n X_i$ is unbiased. Here $X_{(n)} = \max\{X_1, \dots, X_n\}$ is a minimal sufficient statistics and can be showed to be complete. The Rao-Blackwellized version of this estimator is $\frac{n+1}{n} X_{(n)}$.

Remark 2.6.6. We may have a situation when no unbiased estimator exists; e.g. estimating θ in a $\text{Bin}(n, \frac{1}{\theta})$. Indeed, we would require

$$\theta = \mathbb{E}\delta(X) = \sum_{k=0}^n \delta(k) \binom{n}{k} \frac{1}{\theta^k} (1 - \frac{1}{\theta})^{n-k}.$$

After multiplying by θ^n we get that

$$\theta^{n+1} = \sum_{k=0}^n \delta(k)^n \binom{n}{k} (\theta - 1)^{n-k},$$

which is impossible to hold for all $\theta \in (0, 1)$ irrespective of δ because on the right we have a polynomial of order n .

In the special case when $\theta \in \mathbb{R}$ with the square loss function, if δ is unbiased then $R(\theta, \delta) = \text{var}(\delta(X))$. In this case δ in Theorem 2.6.5 is called the **unbiased estimator with uniformly minimum variance (UMVU)**.

Example 2.6.7 (Finding UMVU). Let X_1, \dots, X_n be i.i.d. $\text{Exp}(\mu, 1)$. So

$$p(x; \mu) = \begin{cases} \exp(-(x - \mu)) & x \geq \mu \\ 0 & x < \mu. \end{cases}$$

Both $T_1 = X_{(1)} - \frac{1}{n}$ and $T_2 = \frac{1}{n} \sum_i X_i - 1$ are possible estimators and both are unbiased. Moreover, $\text{var}(T_1) = \frac{1}{n^2} \ll \text{var}(T_2) = \frac{1}{n}$. Note that T_1 is a function of a minimal sufficient statistics. To show that $T = X_{(1)}$ is a sufficient statistics we use the Fisher-Neyman factorization theorem and note that the distribution of the data satisfies

$$\prod_{i=1}^n p(x; \mu) = e^{-\sum_i x_i} e^{n\mu} \mathbf{1}\{T \geq \mu\}$$

and so indeed T is sufficient. To show that T is complete, note that the density of T is

$$p_T(t; \mu) = \begin{cases} n \exp(n(\mu - t)) & t \geq \mu \\ 0 & t < \mu. \end{cases}$$

Let $g(T)$ be a measurable function such that $\mathbb{E}_\mu g(T) = 0$ for all $\mu > 0$. Equivalently, for all $\mu > 0$

$$G(\mu) = e^{n\mu} \int_\mu^\infty e^{-nt} g(t) dt = 0.$$

In particular, G is almost everywhere differentiable with $G'(\mu) = 0$. By the fundamental theorem of calculus,

$$0 = G'(\mu) = ne^{n\mu} \int_\mu^\infty e^{-nt} g(t) dt + e^{n\mu} (-e^{-n\mu} g(\mu)) = g(\mu)$$

implying that T must be complete. Thus, T_1 is actually the UMVU.

2.6.2 Equivariance constraints*

Equivariance is a classical topic in statistics. Recently it also got a lot of attention in machine learning as equivariant convolutional networks grew popular. The idea is very simple and we will explain it first on a concrete example. Suppose the exercise is to classify pictures into one of the categories. Pictures may get rotated and we do not want the label to depend on this rotation. In other words, we want the procedure to be **invariant** under picture rotations. Suppose now that before running the classification exercise, we first reduce the quality of the images in order to save space. In that case, we may want to require that the compression algorithm on a rotated image outputs rotation of the compression of the original image. In that case, we say that the procedure is **equivariant**.

In order to generalize this simple example, we formalize the setup. Recall that a **group** is a set G with a binary operation “ \cdot ” and identity element $e \in G$ such that:

- (i) $g \cdot h \in G$ for all $g, h \in G$,
- (ii) $(g \cdot h) \cdot k = g \cdot (h \cdot k)$ for all $g, h, k \in G$. (associativity)
- (iii) $g \cdot e = e \cdot g = g$ for all $g \in G$,
- (iv) for every $g \in G$ there exists $h \in G$ such that $g \cdot h = h \cdot g = e$, we write $h = g^{-1}$.

The following examples will be important in the sequel.

Example 2.6.8. The set of real numbers \mathbb{R} with addition forms a group. Its identity element is 0. We denote this group by $(\mathbb{R}, +)$.

Example 2.6.9. The group of permutations of the set $\{1, \dots, m\}$.

Example 2.6.10. The set of real invertible $m \times m$ matrices forms a group under the matrix multiplication with the identity element given by the identity matrix. This group is sometimes called the generalized linear group denoted $GL_m(\mathbb{R})$.

Example 2.6.11. The group $SO(m)$ of rotations in \mathbb{R}^m . The group is a subgroup of $GL_m(\mathbb{R})$.

The groups in this section will act on \mathcal{X} , that is, each $g \in G$ defines a function $g : \mathcal{X} \rightarrow \mathcal{X}$, which (perhaps abusing notation a bit) we denote in the same way, and:

1. $e(x) = x$ for all $x \in \mathcal{X}$,
2. $(gh)(x) = g(h(x))$ for all $g, h \in G$ and $x \in \mathcal{X}$.

Following the standard algebraic notation, we write $g \cdot x$ for this action. For example, the group $GL_m(\mathbb{R})$ acts on \mathbb{R}^m by the matrix multiplication $\mathbf{x} \mapsto g \cdot \mathbf{x}$ for a matrix $g \in GL_m(\mathbb{R})$. The group $(\mathbb{R}, +)$ also acts on \mathbb{R}^m by translations $\mathbf{x} \mapsto \mathbf{x} + c\mathbf{1}$ for $c \in (\mathbb{R}, +)$. The group of permutations acts on \mathbb{R}^m by permuting the coordinates. If σ is a permutation of $\{1, \dots, m\}$ then the corresponding transformation of \mathbb{R}^m is $\mathbf{x} \mapsto (x_{\sigma(1)}, \dots, x_{\sigma(m)})$.

The same abstract group can act on different sets. Fix a sample space $\mathcal{X} \subseteq \mathbb{R}^m$ and a group G that acts on it. Suppose $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ is a model for a random variable $X \in \mathcal{X}$ with the identifiable parameter θ (i.e. $\mathbb{P}_\theta = \mathbb{P}_{\theta'}$ implies $\theta = \theta'$). We say that \mathcal{P} is a **group transformation model with respect to G** if it holds that:

$$X \sim \mathbb{P}_\theta \quad \text{then} \quad gX \sim P_{\theta'} \quad \text{for some } \theta' \in \Theta.$$

(the model is invariant under the group action) This particular θ' then determined by θ and the transformation g . In other words, the transformations g also act on θ . We will use the same notation for the action on \mathcal{X} and on Θ :

$$X \sim \mathbb{P}_\theta \quad \iff \quad g \cdot X \sim P_{g \cdot \theta}.$$

Example 2.6.12. Let \mathbb{P}_0 be a probability measure with symmetric density p_0 with respect to the Lebesgue measure on \mathbb{R} (the mean of this distribution is zero if it exists). For $X \sim \mathbb{P}_0$ and $\theta \in \mathbb{R}$, let \mathbb{P}_θ be the distribution of $X + \theta$. Doing this for all θ generates the family $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \mathbb{R}\}$. The parameter θ is called the **location parameter**. The normal family $N(\theta, 1)$ is a special case. Here the group is $(\mathbb{R}, +)$ and we have

$$X \sim \mathbb{P}_\theta \quad \iff \quad X + c \sim P_{\theta+c}.$$

Example 2.6.13. Extending the previous example, suppose that $\mathbb{P}_{0,1}$ is a distribution of a univariate random variable X with mean zero and variance 1. Denote by $\mathbb{P}_{\mu,\sigma}$ the distribution of $\sigma X + \mu$ for $\mu \in \mathbb{R}$ and $\sigma > 0$. Then (μ, σ) is called the **location-scale parameter** for X . This setup is generalized to a random vector X in \mathbb{R}^m by considering the distributions of $\sigma X + \mu \mathbf{1}$. In the vector case it is more suitable to consider a location-scale family by taking $\boldsymbol{\mu} \in \mathbb{R}^m$ and $U \in \mathbb{S}_+^m$ and considering the induced distributions of $UX + \boldsymbol{\mu}$. What is the associated group?

Example 2.6.14. Suppose X has a m -variate zero-mean Gaussian distribution with covariance matrix Σ . The group $\text{GL}_m(\mathbb{R})$ or $m \times m$ invertible matrices acts on $\mathcal{X} = \mathbb{R}^m$, $\mathbf{x} \mapsto A\mathbf{x}$, $A \in \text{GL}_m(\mathbb{R})$. If X is Gaussian with covariance Σ , then AX is Gaussian with covariance $A\Sigma A^\top$ and so the action of $\text{GL}_m(\mathbb{R})$ on $\Theta = \mathbb{S}_+^m$ is $A \cdot \Sigma = A\Sigma A^\top$. We have

$$X \sim \mathbb{P}_\Sigma \quad \Leftrightarrow \quad AX \sim \mathbb{P}_{A\Sigma A^\top}.$$

Example 2.6.15. A slightly more involved example considers a simple Gaussian graphical model given by three-dimensional centered Gaussian distributions with the inverse covariance matrix satisfying $K_{13} = 0$. The group is given by all invertible matrices of the form $g \cdot h$, where

$$h \in \left\{ \left[\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right], \left[\begin{array}{ccc} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right] \right\}, \quad g = \begin{bmatrix} * & * & 0 \\ 0 & * & 0 \\ 0 & * & * \end{bmatrix}.$$

The action on the parameter space is the same as in the previous example. (verify the details)

Consider a decision procedure in the situation when there is a group acting on \mathcal{X} and Θ as described above. As it was noted in the introductory example, if we have a group acting on the sample space, we may want some statistical procedures like testing or classification to be invariant with respect to this group action. A decision rule δ is called **invariant** if $\delta(g \cdot x) = \delta(x)$ for all $x \in \mathcal{X}$ and all $g \in G$.

Example 2.6.16. Suppose $X \sim \mathbb{P}_\theta$ and the goal is to test hypothesis $\theta \in \Theta_0$ versus $\theta \in \Theta$, $\mathcal{A} = \{0, 1\}$. Suppose that a group G acts on \mathcal{X} and on Θ_0 , that is, $X \sim \mathbb{P}_\theta$ for $\theta \in \Theta_0$ then $g \cdot X \sim \mathbb{P}_{\theta'}$ for $\theta' \in \Theta_0$. Intuitively, if for data X we accept/reject the null, we should also accept/reject it for data $g \cdot X$. We will be then interest in test procedures $\delta : \mathcal{X} \rightarrow \{0, 1\}$ that are invariant.

Example 2.6.17. Consider a classification function $\delta : \mathcal{X} \rightarrow \{1, \dots, c\}$. In many applications, there is a group of transformations of \mathcal{X} and we want the classification procedure to classify each example the same way as its transformed version (think images and their rotated versions). We again will require that δ is invariant.

In estimation of group transformation models we have $\mathcal{A} = \Theta$ and in this case G acts on \mathcal{A} in the same way as it acts on Θ . In this case we may want to restrict to procedures that are equivariant.

Definition 2.6.18. A function $\delta : \mathcal{X} \rightarrow \Theta$ is *equivariant* if $\delta(g \cdot x) = g \cdot \delta(x)$.

Example 2.6.19. Let $\mathcal{X} = \mathbb{R}^m$ and define maps $g_c(\mathbf{x}) = \mathbf{x} + c\mathbf{1}$, the location shifts. These transformations form a group. The function $\beta(\mathbf{x}) = \bar{\mathbf{x}}\mathbf{1}$ is equivariant where $\bar{\mathbf{x}}$ is the average of entries of \mathbf{x} . The function $\alpha(\mathbf{x}) = \mathbf{x} - \bar{\mathbf{x}}\mathbf{1}$ is invariant.

Although this plays no role here, we note that equivariance is also a convex constraint under a minor condition.

← Exercise 2.7.18

The last ingredient in designing good equivariant procedures is a loss function that is also amenable to this group action setting.

Definition 2.6.20. The loss function $L : \Theta \times \mathcal{A} \rightarrow [0, \infty]$ is called *invariant* (with respect to G) if, for each $g \in G$, $\theta \in \Theta$, and $a \in \mathcal{A}$, $L(g \cdot \theta, g \cdot a) = L(\theta, a)$.

Example 2.6.21. [Example 2.6.14 continued] Let

$$L(\Sigma, S) = -\log \det(S\Sigma^{-1}) + \text{tr}(S\Sigma^{-1} - I_m)$$

be the Kullback-Leibler divergence between two mean-zero Gaussian distributions with covariances S and Σ . Basic matrix algebra gives that

$$L(g\Sigma g^\top, gSg^\top) = L(\Sigma, S).$$

Thus, in the problem of estimating Σ we have $\mathcal{A} = \mathbf{S}_+^m$ and $L(\Sigma, S)$ is a valid invariant loss function.

← Exercise 2.7.19

Definition 2.6.22. An *invariant decision problem* is when: \mathcal{P} is a group transformation model, δ is equivariant, and L is invariant.

We view the insistence that the decision rule be equivariant as a constraint on the possible decision rules, just like unbiasedness is a constraint. Then the question is if there is an equivariant rule that uniformly minimizes the risk. The first result is a step in this direction.

Theorem 2.6.23. In an invariant decision problem, the risk function $R(\theta, \delta)$ of an equivariant decision rule δ is an invariant function on Θ .

Proof. We have

$$\begin{aligned} R(g \cdot \theta, \delta) &= \mathbb{E}_{g \cdot \theta}(L(g \cdot \theta, \delta(X))) = \mathbb{E}_{g \cdot \theta}(L(\theta, g^{-1} \cdot \delta(X))) \\ &= \mathbb{E}_{g \cdot \theta}(L(\theta, \delta(g^{-1} \cdot X))) = \mathbb{E}_\theta(L(\theta, \delta(X))), \end{aligned}$$

where the last equation follows from the fact that $X \sim \mathbb{P}_{g \cdot \theta}$ if and only if $g^{-1} \cdot X \sim \mathbb{P}_\theta$. □

This theorem really shows that $R(\theta, \delta)$ is constant on G -orbits in Θ , that is, the sets

$$G \cdot \theta = \{\theta' \in \Theta : \theta' = g \cdot \theta \text{ for some } g \in G\}.$$

This gives a plausible justification for restricting attention to equivariant decision rules. Since the risk function is constant on orbits in Θ when the loss is invariant, this makes it easier to compare equivariant rules by means of their risk functions. In particular, if the group acts transitively on the parameter space (i.e. there is only one orbit), then the problem of noncomparability of risk functions disappears altogether as the risk function becomes constant on δ . This happens in some interesting situations, e.g. for the location model.

Theorem 2.6.24 (Pitman's estimator). *Suppose that $Y = (X_1 - X_n, \dots, X_{n-1} - X_n)$ and $L(\theta, a) = (\theta - a)^2$. Suppose that δ_0 is a location equivariant estimator with finite risk. Then, the equivariant estimator with smallest risk is $\delta_0(X) - \mathbb{E}_0[\delta_0(X)|Y]$.*

Proof. If δ_0 is a location equivariant estimator with finite risk then all other equivariant estimators have the form $\delta_0(X) - v(Y)$. Indeed, δ is equivariant if and only if $\delta - \delta_0$ is invariant. So it remains to show that every invariant function f is a function of Y . This follows because

$$f(\mathbf{x}) = f((\mathbf{y}, 0) + x_n \mathbf{1}) = f(\mathbf{y}, 0) = v(\mathbf{y}).$$

By Theorem 2.6.23, the risk function is invariant in θ for an equivariant δ , and so

$$R(\theta, \delta) = R(0, \delta) = \mathbb{E}_0[\delta_0(X) - v(Y)]^2 = \mathbb{E}_0(\mathbb{E}_0[(\delta_0(X) - v(Y))^2|Y]),$$

which is minimized by minimizing $\mathbb{E}_0[(\delta_0(X) - v(Y))^2|Y = \mathbf{y}]$ uniformly in \mathbf{y} . This is accomplished by choosing $v(\mathbf{y}) = \mathbb{E}_0[\delta_0(X)|Y = \mathbf{y}]$. \square

Like with the Lehmann-Scheffe theorem, we get an explicit procedure of obtaining the best equivariant estimator by improving any given equivariant estimator.

It can be shown that the Pitman's estimator is the generalized Bayes rule with respect to the uniform "prior" distribution. This fact and Theorem 2.6.24 can be both generalized; see Section 6.2.3 in ⁸; we will not provide any more details.

Theorem 2.6.25. *Consider an invariant decision problem. Under some assumptions, if the formal Bayes rule with respect to the right invariant Haar prior on G exists, then it is the minimum risk equivariant rule.*

The following case of the truncated normal distribution shows the power of this result. Using Theorem 2.6.24 directly is rather tedious. Computing the generalized Bayes rule $\mathbb{E}(\theta|X)$ in this case is straightforward.

← Exercise 2.7.20

⁸ Mark J. Schervish. *Theory of statistics*. Springer Series in Statistics. Springer-Verlag, New York, 1995

← Exercise 2.7.21

2.6.3 Type I error constraints

Recall the testing problem in Example 2.2.2, which is a particular instance of a decision problem with $\mathcal{A} = \{0, 1\}$ and 0/1-loss

$$L(\theta, a) = \mathbb{1}(a = 1, \theta \in \Theta_0) + \mathbb{1}(a = 0, \theta \in \Theta_1) \quad (2.15)$$

which is associated to testing two competing hypotheses $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$. The data are $X \sim \mathbb{P}_\theta$ for $\theta \in \Theta_0 \cup \Theta_1$.

It is usual to find tests that bound the **type I error**, defined as $\sup_{\theta \in \Theta_0} \beta(\theta)$, on some level α , where the power function $\beta(\cdot)$ is defined in (2.3). By the risk calculation in Example 2.2.2,

$$\sup_{\theta \in \Theta_0} \beta(\theta) = \sup_{\theta \in \Theta_0} R(\theta, \delta).$$

On the other hand

$$\sup_{\theta \in \Theta_1} R(\theta, \delta) = \sup_{\theta \in \Theta_1} (1 - \beta(\theta)) = 1 - \inf_{\theta \in \Theta_1} \beta(\theta),$$

which is precisely the **type II error**. It follows that

$$\bar{R}(\delta) = \sup_{\theta \in \Theta_0 \cup \Theta_1} R(\theta, \delta) = \max\{\sup_{\theta \in \Theta_0} R(\theta, \delta), \sup_{\theta \in \Theta_1} R(\theta, \delta)\}.$$

By “bounding the type I error” we mean optimizing $\bar{R}(\delta)$ by restricting to test procedures δ with an explicit constraint on the type I error:

$$\sup_{\theta \in \Theta_0} R(\theta, \delta) \leq \alpha.$$

← Exercise 2.7.22

Definition 2.6.26. A test φ^* with level α is called **uniformly most powerful (UMP)** if

$$\mathbb{E}_\theta \varphi^* \geq \mathbb{E}_\theta \varphi, \quad \text{for all } \theta \in \Theta_1$$

for all φ with level at most α .

As we will see, test of this form may appear only in very special situations and mostly in the univariate case. A particularly famous instance is discussed in Section 3.1.

2.7 Exercises

Exercise 2.7.1. Show that if $f : \mathcal{A} \rightarrow [-\infty, \infty]$ for $\mathcal{A} \subseteq \mathbb{R}^d$ is lower semicontinuous then it is measurable. Hint: Use Exercise A.2.8.

Exercise 2.7.2. Suppose that $x_1, \dots, x_n \in \mathbb{R}$. Show that the minimizer of $f(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta - x_i)^2$ is the average $\frac{1}{n} \sum_{i=1}^n x_i$. Moreover, show that the median of the collection x_1, \dots, x_n minimizes $f(\theta) = \frac{1}{n} \sum_{i=1}^n |\theta - x_i|$.

Exercise 2.7.3. Show that the estimator δ_3 in Example 2.1.1 is a Bayes rule for the quadratic loss and a beta distribution.

The next two exercises allow you to explore a similar situation in the case when the loss function is not differentiable.

Exercise 2.7.4. In the one-dimensional case, show that if $L(\theta, a) = |\theta - a|$ then the Bayes estimator is the median of the posterior distribution $\pi(\theta|x)$ (c.f. Exercise 2.7.2).

Exercise 2.7.5. Show that if $L(\theta, a)$ is the 0/1-loss then the Bayes estimator is the mode of the posterior distribution $\pi(\theta|x)$.

Exercise 2.7.6. In the setting of Theorem 2.3.7 show that π is least favourable.

Exercise 2.7.7. Show that in Section 2.3.5 we always have $L^* \leq U^*$.

Exercise 2.7.8. Let (X_1, \dots, X_n) be a random sample of binary random variables with $X_i \sim \text{Bern}(\theta)$ with $\theta \in (0, 1)$.

(i) Show that the sample mean \bar{X}_n is an admissible estimator of θ under the loss function $L(\theta, a) = (a - \theta)^2 / [\theta(1 - \theta)]$.

(ii) Show that \bar{X} is an admissible estimator of θ under the squared error loss.

Exercise 2.7.9. Show that if a minimax rule is essentially unique then it is admissible.

Exercise 2.7.10. Show the following result: If risk functions for all decision rules are continuous in θ , if δ is Bayes for π and has finite integrated risk $r(\pi, \delta) < +\infty$, and if the support of π is the whole parameter space, then δ is admissible. Hint: This looks much more complicated than it actually is.

Exercise 2.7.11. Show that if $\delta(X)$ is an unbiased estimator, $\mathbb{E}_\theta \delta(X) = \theta$, then $\eta(T)$ is also unbiased.

Exercise 2.7.12. Consider the special case when $\theta = 0$. In this case $\|X\|^2 \sim \chi_d^2$. Show that $\mathbb{E}_0 \frac{1}{\|X\|^2} = \frac{1}{d-2}$ and so

$$R(0, \delta_{\text{JS}}) = 2 \ll d.$$

Exercise 2.7.13. Suppose that $\mathbf{Z} \in \mathbb{R}^{d \times p}$ is a fixed matrix of full column rank (in particular $p \leq d$). Consider a linear estimator of μ given by $\hat{\mu} = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top X$. Compute the risk of this estimator. Discuss conditions when the risk of this estimator can be significantly smaller given $p \ll d$.

Exercise 2.7.14. Show that the linear estimator $\hat{\mu}_C$ with C defined in (2.14) satisfies the conditions of Proposition 2.5.6 and so it is admissible.

Exercise 2.7.15. Show that the decision rules that are unbiased with respect to the square loss $L(\theta, \delta) = \|\theta - \delta(X)\|^2$ are the rules satisfying $\mathbb{E}_\theta \delta(X) = \theta$.

Exercise 2.7.16. Suppose $Y \sim N_n(\mu, I_n)$. Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ be a fixed matrix with full column rank. In this exercise we consider estimators of the form $\hat{\mu} = \mathbf{X}\hat{\beta}$ for some estimator $\hat{\beta}$. More specifically, we study the least squares estimator $\hat{\beta}^{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top Y$.

- (i) Find the formula for the risk $R(\mu, \hat{\mu}_C)$ of a general linear estimator $\hat{\mu}_C = CY$.
- (ii) Let $\hat{\mu}^{\text{LS}} = \mathbf{X}\hat{\beta}^{\text{LS}}$ be the corresponding estimator of the mean of Y . Find the risk of this estimator directly in terms of $C_0 = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, μ , n , and p only.
- (iii) Suppose $\mu = \mathbf{X}\beta^*$ for some $\beta^* \in \mathbb{R}^p$. Show that $R(\mu, \hat{\mu}^{\text{LS}}) = p$.
- (iv) Consider the ridge estimator $\hat{\mu}^{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \delta I_p)^{-1} \mathbf{X}^\top Y$ and the corresponding matrix $C_\delta = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \delta I_p)^{-1} \mathbf{X}^\top$, where $\delta \geq 0$. Show that $R(\mu, \hat{\mu}^{\text{ridge}}) = \text{tr}(C_\delta^2) + \|(I_n - C_\delta)\mu\|^2$.
- (v) Suppose that $\mathbf{X}^\top \mathbf{X} = I_p$ and that $\mu = \mathbf{X}\beta^*$ for some β^* . Show that for every $\beta^* \neq 0$ there exists $\delta > 0$ such that $R(\mathbf{X}\beta^*, \hat{\mu}^{\text{ridge}}) < R(\mathbf{X}\beta^*, \hat{\mu}^{\text{LS}})$.

Exercise 2.7.17. Consider estimation in a regular exponential family and let the action space \mathcal{A} be its set of mean parameters. Consider the loss function given by the Kullback-Leibler divergence in (1.18):

$$L(\theta_2, \mu_1) = K(\mu_1, \theta_2) = A^*(\mu_1) + A(\theta_2) - \langle \theta_2, \mu_1 \rangle.$$

Given a sample X_1, \dots, X_n , consider the sample average of the sufficient statistics $\delta(X_1, \dots, X_n) = \bar{\mu}_n$. The minimizer $\hat{\theta}_n$ of $L(\theta, \bar{\mu}_n)$ is the MLE of the canonical parameter. Although $\bar{\mu}_n$ is an unbiased estimator (in the usual sense) of the true mean parameter μ , this is typically not true for $\hat{\theta}_n$. Show that $\bar{\mu}_n$ is an unbiased decision rule in the sense of Definition 2.6.1.

Exercise 2.7.18. Show that the set of equivariant randomized decision rules is convex (with convex combination defined as in (2.4)). Discuss some natural conditions under which the set of non-randomized equivariant procedures also forms a convex set.

Exercise 2.7.19. Show that the squared loss function is invariant for the translation family in Example 2.6.12.

Exercise 2.7.20. Show that for the Gaussian model $N(\theta, 1)$ the Pitman's estimator is the sample mean.

Exercise 2.7.21. Let (X_1, \dots, X_n) be a random sample of random variables with the Lebesgue density

$$f_\theta(x) = \begin{cases} \sqrt{\frac{2}{\pi}} e^{-(x-\theta)^2/2} & \text{if } x \geq \theta, \\ 0 & \text{otherwise} \end{cases}$$

where $\theta \in \mathbb{R}$ is unknown. Find the minimum risk location equivariant estimator of θ under the squared loss.

Exercise 2.7.22. Show that bounding the type I error in Section 2.6.3 gives a convex constraint on $\delta \in \mathcal{D}$.

3

Hypothesis testing and multiple testing (2 weeks)

Our goal in this section is not to give an extensive treatment of hypothesis testing procedures. We trust that much of it was covered in earlier courses. This includes the basic philosophy of constructing statistical tests (rejection regions, p-values etc).

Given data $X \sim \mathbb{P}_\theta$ and a model $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$, consider two competing hypothesis: $H_0 : \theta \in \Theta_0, H_1 : \theta \in \Theta_1$. Unless otherwise stated we assume that $\Theta = \Theta_0 \cup \Theta_1$ and $\Theta_0 \cap \Theta_1 = \emptyset$. In this context, $\mathcal{A} = \{0, 1\}$ and we are interested in statistical testing procedures $\delta : \mathcal{X} \rightarrow \{0, 1\}$. In consequence, every **nonrandomized test** can be associated with a measurable set $S \subseteq \mathcal{X}$ such that the decision rule becomes: accept H_1 if $X \in S$ and accept H_0 when $X \notin S$. From the theoretical perspective it is also useful to allow for **randomized tests**, in which case the **critical function** φ to every $x \in \mathcal{X}$ is assigns the probability of rejecting H_0 given $X = x$

$$\varphi(x) = \mathbb{P}(\text{reject } H_0 | X = x),$$

where φ is assumed to be a measurable function from \mathcal{X} to $[0, 1]$. A nonrandomized test is a randomized test with $\varphi(X) = \mathbb{1}(X \in S)$.

Recall from Example 2.2.2 that the **power function** is the function $\beta : \Theta \rightarrow [0, 1]$ defined by

$$\beta(\theta) := \mathbb{P}_\theta(\text{reject } H_0) = \int \varphi(x) d\mathbb{P}_\theta = \mathbb{E}_\theta \varphi(X).$$

Using the 0/1-loss we see that the risk function satisfies

$$R(\theta, \delta) = \begin{cases} \beta(\theta) & \text{if } \theta \in \Theta_0 \\ 1 - \beta(\theta) & \text{if } \theta \in \Theta_1 \end{cases}.$$

In the idealised situation we would have $\beta(\theta) = 0$ if $\theta \in \Theta_0$ and $\beta(\theta) = 1$ for $\theta \in \Theta_1$. This is obviously not possible in most of the cases. For example, when all \mathbb{P}_θ have the same support. However, a good test would try to be close to that. We define the **significance**

level of a test as

$$\alpha = \sup_{\theta \in \Theta_0} \beta(\theta) = \sup_{\theta \in \Theta_0} \mathbb{E}_\theta \varphi(X).$$

The main technical convenience of randomized tests comes from the fact that the set of critical functions is a convex set.

3.1 Neyman-Pearson lemma*

In this section we discuss a classical result in the case when H_0 and H_1 are both simple with the underlying density functions p_0 and p_1 with respect to some underlying measure μ . The type I error is

$$\mathbb{E}_0 \varphi(X) = \int \varphi(x) p_0(x) d\mu(x)$$

and another quantity of interest is

$$\mathbb{E}_1 \varphi(X) = \int \varphi(x) p_1(x) d\mu(x).$$

A good test will result in small type I error (small $\mathbb{E}_0 \varphi$) and small type II error (big $\mathbb{E}_1 \varphi$). We want to explore to what extent this can be achieved.

As discussed in Section 2.6.3, a natural approach is to consider the optimization problem:

$$\text{maximize } \mathbb{E}_1 \varphi(X) \quad \text{subject to } \mathbb{E}_0 \varphi(X) \leq \alpha.$$

We call such test the most powerful.

Consider the **likelihood ratio statistic**

$$\text{LR}(x) = \frac{p_1(x)}{p_0(x)},$$

where we define $\text{LR}(x) = +\infty$ if $p_1(x) > 0$, $p_0(x) = 0$, and $\text{LR}(x) = 0$ if $p_0(x) = p_1(x) = 0$. For any $\lambda \geq 0$ we consider the critical function

$$\varphi_\lambda(x) := \begin{cases} 1 & \text{if } \text{LR}(x) > \lambda \\ \gamma & \text{if } \text{LR}(x) = \lambda \\ 0 & \text{if } \text{LR}(x) < \lambda. \end{cases} \quad (3.1)$$

Note that, if $\gamma \in (0, 1)$, the corresponding test is randomized, as we need to toss a γ -coin when we observe $\text{LR}(x) = \lambda$. Our freedom to choose γ will add extra flexibility in a second.

Constructing a test of size α is not a problem. It is enough to toss an α -coin to decide on the decision. However, this procedure does not depend on the data so it is clear that its power will be low.

Theorem 3.1.1 (Neyman-Pearson Lemma). *Given any level $\alpha \in [0, 1]$ there exists a likelihood ratio test φ with level α . Any likelihood ratio test with level α maximizes $\mathbb{E}_1\varphi$ among all tests with level $\leq \alpha$.*

Proof. To prove the first part note that for $\alpha = 0$ we can take $\lambda = +\infty$, and for $\alpha = 1$ we can take $\gamma = 1, \lambda = 0$ (check carefully). Now take $\alpha \in (0, 1)$. By construction, $\mathbb{P}_0(\text{LR} = \infty) = 0$ and so LR is finite \mathbb{P}_0 -a.s. We claim that this implies that

$$\exists \lambda < \infty \text{ such that } \mathbb{P}_0(\text{LR} > \lambda) \leq \alpha \text{ and } \mathbb{P}_0(\text{LR} \geq \lambda) \geq \alpha. \tag{3.2}$$

Before we continue, take a look at Figure 3.1. As we showed in Proposition C.1.2, the survival function $G(\lambda) = \mathbb{P}_0(\text{LR} > \lambda)$ is right-continuous. Thus taking

$$\lambda := \inf_{t \geq 0} \{t : \mathbb{P}_0(\text{LR} > t) \leq \alpha\}$$

we get $\mathbb{P}_0(\text{LR} > \lambda) \leq \alpha$. By Remark C.1.3, the function $\mathbb{P}_0(\text{LR} \geq \lambda)$ is left-continuous in λ . Since, for every $\epsilon > 0$

$$\mathbb{P}_0(\text{LR} \geq \lambda - \epsilon) \geq \mathbb{P}_0(\text{LR} > \lambda - \epsilon) > \alpha,$$

we get that $\lim_{t \rightarrow \lambda^-} \mathbb{P}_0(\text{LR} \geq t) = \mathbb{P}_0(\text{LR} \geq \lambda) \geq \alpha$ confirming (3.2).

By properly randomizing our procedure, we can obtain level exactly α . If $\mathbb{P}_0(\text{LR} = \lambda) = 0$ then $\mathbb{P}_0(\text{LR} > \lambda) = \mathbb{P}_0(\text{LR} \geq \lambda) = \alpha$ and so the likelihood ratio test for this λ has level exactly α . If $\mathbb{P}_0(\text{LR} = \lambda) > 0$, we define

$$\gamma = \frac{\alpha - \mathbb{P}_0(\text{LR} > \lambda)}{\mathbb{P}_0(\text{LR} \geq \lambda) - \mathbb{P}_0(\text{LR} > \lambda)} \in (0, 1)$$

and it is straightforward to check that $\mathbb{E}_0\varphi_\lambda(X) = \alpha$ with this choice of γ .

By the first part, there exists $\lambda \geq 0$ (and $\gamma \in [0, 1]$) such that φ_λ has size α . Let φ be any test with level $\leq \alpha$. We have

$$\begin{aligned} \mathbb{E}_1\varphi(X) &\leq \mathbb{E}_1\varphi(X) - \lambda(\mathbb{E}_0\varphi - \alpha) = \int \varphi(x)(p_1(x) - \lambda p_0(x))\mu(dx) + \lambda\alpha \\ &\leq \int \varphi_\lambda(x)(p_1(x) - \lambda p_0(x))\mu(dx) + \lambda\alpha \\ &\leq \mathbb{E}_1\varphi_\lambda(X) - \lambda(\mathbb{E}_0\varphi_\lambda - \alpha) = \mathbb{E}_1\varphi_\lambda(X). \end{aligned}$$

□

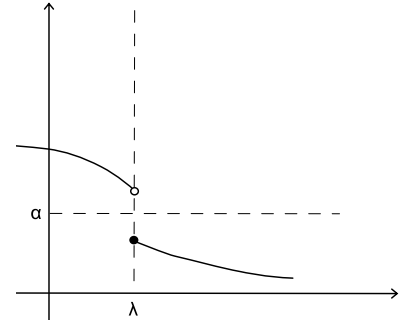


Figure 3.1: The claim (3.2) is clear if the CDF $\mathbb{P}(\text{LR} \leq \lambda)$ is continuous. If it is not and α occurs at one of the jumps, we use the fact that the CDF, and so also the function $\mathbb{P}_0(\text{LR} > \lambda)$ is right-continuous as on the picture.

← Exercise 3.7.1

The following corollary will be useful.

Lemma 3.1.2. *In the likelihood ratio test with level $\alpha = \mathbb{E}_0\varphi_\lambda(X)$ we have $\mathbb{E}_1\varphi_\lambda(X) > \alpha$ with equality if and only if $p_0 \equiv p_1$.*

Proof. A silly test δ , where we decide based on an α -coin flip regardless the data ($\varphi(x) = \alpha$ for all x) satisfies $\alpha = \mathbb{E}_0\delta(X)$ but also $\mathbb{E}_1\delta(X) = \alpha$. By the Neyman-Pearson lemma we conclude that $\mathbb{E}_1\varphi_\lambda(X) \geq \alpha$. By Lemma 3.7.1, we can conclude that the inequality is strict, unless the silly coin-flip is equivalent with the likelihood ratio test, which happens if and only if $p_0(x) = p_1(x)$ for all x . \square

Example 3.1.3. Suppose $p_\theta(x) = \theta e^{-\theta x}$ for $x > 0$. Consider the test $H_0 : \theta = 1$ versus $H_1 : \theta = \theta_1$ for a fixed $\theta_1 > 1$. A likelihood ratio test is given by the rejection region

$$\text{LR}(x) = \frac{p_1(x)}{p_0(x)} = \theta_1 e^{(1-\theta_1)x} > \lambda$$

or equivalently if

$$x < \frac{\log(\theta_1/\lambda)}{\theta_1 - 1} = \lambda'.$$

The level is

$$\alpha = \mathbb{P}_0(X < \lambda') = \int_0^{\lambda'} e^{-x} dx = 1 - e^{-\lambda'}.$$

Solving $\lambda' = -\log(1 - \alpha)$ gives a test with level α . By Proposition 3.1.1 the test constructed in this way maximizes $\mathbb{E}_1\varphi$ over all tests with level at most α . Note that the test does not depend on θ_1 ! Hence, this test is the UMP test for $H_0 : \theta = 1$ versus $H_1 : \theta_1 > 1$. In the next section we provide some general theory explaining this phenomenon.

Often, instead of $\text{LR}(x)$ we work with the log-likelihood ratio $\lambda(x) = \log \text{LR}(x)$ (for essentially the same reason as we prefer log-likelihoods). If p_0, p_1 both lie in the same exponential family with parameters θ_0, θ_1 then

$$\lambda(\mathbf{x}) = \langle \theta_1 - \theta_0, \mathbf{t}(\mathbf{x}) \rangle - (A(\theta_1) - A(\theta_0))$$

and so thresholding $\text{LR}(\mathbf{x})$ is equivalent to thresholding a linear function of $\mathbf{t}(\mathbf{x})$. The following simple example plays an important role later.

Example 3.1.4 (Linear discriminant analysis). Consider two classes that are distributed as multivariate Gaussians, say $N(\mu_0, \Sigma)$ and $N(\mu_1, \Sigma)$, respectively, differing only in their mean vectors. For each observation we want to decide from which of the two classes it comes. In this case, the log-likelihood ratio reduces to the linear statistic

$$\lambda(\mathbf{x}) = \left\langle \mu_1 - \mu_0, \Sigma^{-1} \left(x - \frac{\mu_0 + \mu_1}{2} \right) \right\rangle. \quad (3.3)$$

(Note that $\lambda(\mu_1) > 0$ and $\lambda(\mu_0) < 0$)

By the Neyman-Pearson lemma, the optimal decision rule is based on thresholding this statistic. Concretely, if the two classes are equally likely, then the corresponding Bayes risk of our procedure is given by

$$\text{Err}(\lambda) := \frac{1}{2}\mathbb{P}_0(\lambda(X) \geq \lambda) + \frac{1}{2}\mathbb{P}_1(\lambda(X) \leq \lambda) \quad (3.4)$$

and we use it to evaluate the quality of this decision rule.

← Exercise 3.7.2

Given our Gaussian assumptions, some algebra shows that the error probability can be written in terms of the Gaussian cumulative distribution function Φ as

$$\text{Err}(0) = \Phi(-\gamma/2) \quad \text{where } \gamma = \|\mu_1 - \mu_0\|_{\Sigma}, \quad (3.5)$$

where $\|x\|_{\Sigma} = \sqrt{x^{\top}\Sigma^{-1}x}$.

3.1.1 Derivation from the first principles*

In our calculations above we used a guess that the optimal test should be based on the likelihood ratio and then it was straightforward to show that such a test must be optimal. In this subsection, we argue how we could come up with this guess. Note that our optimization problem is equivalent to maximizing

$$f(\varphi) = \begin{cases} \mathbb{E}_1\varphi & \text{if } \mathbb{E}_0\varphi \leq \alpha, \\ -\infty & \text{otherwise.} \end{cases}$$

Denote by φ^* an optimizer of this function and note that

$$\inf_{\lambda \geq 0} \left\{ \mathbb{E}_1\varphi - \lambda(\mathbb{E}_0\varphi - \alpha) \right\} = f(\varphi). \quad (3.6)$$

The function $L(\varphi, \lambda) = \mathbb{E}_1\varphi - \lambda(\mathbb{E}_0\varphi - \alpha)$ is called the Lagrangian and (3.6) shows that

$$\sup_{\varphi} f(\varphi) = \sup_{\varphi} \inf_{\lambda \geq 0} L(\varphi, \lambda), \quad (3.7)$$

where the supremum over φ is unrestricted and it runs over all critical functions. It is clear that

$$\sup_{\varphi} \inf_{\lambda \geq 0} L(\varphi, \lambda) \leq \inf_{\lambda \geq 0} \sup_{\varphi} L(\varphi, \lambda). \quad (3.8)$$

The next result states that the inequality in (3.8) is actually an equality. The proof is similar to the proof of Theorem 2.3.9.

Proposition 3.1.5. *We have*

$$\sup_{\varphi} \inf_{\lambda \geq 0} L(\varphi, \lambda) = \inf_{\lambda \geq 0} \sup_{\varphi} L(\varphi, \lambda). \quad (3.9)$$

Proof. This is a standard strategy in proving strong duality. Let

$$\mathcal{A} = \{(u, t) : \mathbb{E}_0\varphi - \alpha \leq u, \mathbb{E}_1\varphi \geq t \text{ for some } \varphi\}$$

so that $\mathbb{E}_1\varphi^* = \sup\{t : (0, t) \in \mathcal{A}\}$. Consider also a set

$$\mathcal{B} = \{(0, s) : s > \mathbb{E}_1\varphi^*\}.$$

The sets \mathcal{A}, \mathcal{B} are both convex and disjoint subsets of \mathbb{R}^2 (make sure you agree). By the separating hyperplane theorem (c.f. Theorem B.1.3) there exist real numbers $\tilde{\lambda}, \mu, \beta$ such that

$$\tilde{\lambda}u - \mu t \geq \beta \quad \text{for all } (u, t) \in \mathcal{A}, \quad (3.10)$$

$$\tilde{\lambda}u - \mu t \leq \beta \quad \text{for all } (u, t) \in \mathcal{B}. \quad (3.11)$$

If $(u, t) \in \mathcal{A}$ then $(u', t') \in \mathcal{A}$ for every $u' > u$ and $t' < t$. Hence, (3.10) implies that $\tilde{\lambda} \geq 0$ and $\mu \geq 0$ for otherwise this expression could not be bounded below. Since $u = 0$ in \mathcal{B} , (3.11) states that $-\mu t \leq \beta$ for every $t > \mathbb{E}_1\varphi^*$. But then also $-\mu\mathbb{E}_1\varphi^* \leq \beta$. It follows that for every φ

$$\tilde{\lambda}(\mathbb{E}_0\varphi - \alpha) - \mu\mathbb{E}_1\varphi \geq \beta \geq -\mu\mathbb{E}_1\varphi^*.$$

or equivalently, denoting $\lambda^* = \frac{\tilde{\lambda}}{\mu}$,

$$\lambda^*(\mathbb{E}_0\varphi - \alpha) - \mathbb{E}_1\varphi \geq \frac{\beta}{\mu} \geq -\mathbb{E}_1\varphi^*.$$

This gives that

$$\sup_{\varphi} \inf_{\lambda \geq 0} L(\varphi, \lambda) = \mathbb{E}_1\varphi^* \geq \sup_{\varphi} L(\varphi, \lambda^*) \geq \inf_{\lambda \geq 0} \sup_{\varphi} L(\varphi, \lambda).$$

This proves the reverse of (3.8) and thus the equality in (3.9). \square

By equation (3.7) and Proposition 3.1.5, in order to maximize f we can first optimize $L(\varphi, \lambda)$ over φ and then over λ . For any fixed $\lambda \geq 0$ we have

$$L(\varphi, \lambda) = \mathbb{E}_1\varphi - \lambda(\mathbb{E}_0\varphi - \alpha) = \int_{\mathcal{X}} \varphi(x)(p_1(x) - \lambda p_0(x))dx + \lambda\alpha.$$

As we also argued in the previous section, irrespective of α , it is clear that φ that optimizes $L(\varphi, \lambda)$ satisfies

$$\varphi(x) = \begin{cases} 1 & \text{if } p_1(x) - \lambda p_0(x) > 0 \\ 0 & \text{if } p_1(x) - \lambda p_0(x) < 0 \end{cases}. \quad (3.12)$$

If $p_1(x) - \lambda p_0(x) = 0$, $\varphi(x)$ can take any value. So the optimum φ is given by a likelihood ratio test.

3.2 Some constructions of UMP tests*

If the hypotheses are not both simple the situation is in general much more complicated and a uniformly most powerful test (as defined in Section 2.6.3) may be hard or impossible to obtain. In this section we discuss some of the special situations when a UMP can be obtained. We also provide alternative approaches to find a test with good properties.

3.2.1 Monotone likelihood ratios

A one-dimensional family of densities $p_\theta(x)$, $\theta \in \Theta \subseteq \mathbb{R}$ has **monotone likelihood ratios (MLR)** in $T(x)$ if, whenever $\theta < \theta'$, the likelihood ratio $p_{\theta'}(x)/p_\theta(x) = h(T(x))$ for a nondecreasing function h . A canonical example of such a situation is a one-dimensional exponential family with

$$p_\theta(x) = h(x) \exp\{\eta(\theta)T(x) - A(\theta)\},$$

with the canonical parameter η being a strictly increasing function of the parameter of interest θ . In this case, if $\theta' > \theta$, then

$$\frac{p_{\theta'}(x)}{p_\theta(x)} = \exp\{(\eta(\theta') - \eta(\theta))T(x) + A(\theta) - A(\theta')\},$$

which is increasing in $T(x)$.

In this section we will be interested in testing

$$H_0 : \theta \leq \theta_0 \quad \text{versus} \quad H_1 : \theta > \theta_0.$$

Consider the test

$$\varphi_t(x) = \begin{cases} 1 & \text{if } T(x) > t, \\ \gamma & \text{if } T(x) = t, \\ 0 & \text{if } T(x) < t \end{cases} \quad (3.13)$$

Note that for MLR families the corresponding test is equivalent to the likelihood ratio test. In particular, for any $\theta_1 > \theta_0$, this statistic gives a most powerful test for $H_0 : \theta = \theta_0$ versus $H_1 : \theta = \theta_1$ at level $\alpha = \mathbb{E}_{\theta_0} \varphi_t(X)$. The following result allows us to get a more general statement.

Proposition 3.2.1. *If $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta \subseteq \mathbb{R}\}$ has monotone likelihood ratios and $\mathbb{P}_\theta \neq \mathbb{P}_{\theta'}$ for $\theta \neq \theta'$. Then, for every $t > 0$, the power function $\beta(\theta) = \mathbb{E}_\theta \varphi_t(X)$ is strictly increasing in θ .*

Proof. Note that for any $\theta_1 < \theta_2$, $\varphi_t(X)$ is equivalent to the likelihood ratio test and so it is the most powerful test for testing $H_0 : \theta = \theta_1$ versus $H_1 : \theta = \theta_2$ at level $\mathbb{E}_{\theta_1} \varphi_t(X)$. By Lemma 3.1.2, we conclude that $\mathbb{E}_{\theta_1} \varphi_t(X) \leq \mathbb{E}_{\theta_2} \varphi_t(X)$ with equality if and only if $p_{\theta_1}(x) = p_{\theta_2}(x)$. □

As a corollary from this result we get that

$$\sup_{\theta \leq \theta_0} \mathbb{E}_\theta \varphi_t(X) = \mathbb{E}_{\theta_0} \varphi_t(X).$$

In particular, the same test statistic can be used to get a most powerful test of size $\alpha = \mathbb{E}_{\theta_0} \varphi_t(X)$ for testing $H_0 : \theta \leq \theta_0$ versus $H_1 : \theta = \theta_1$ for any $\theta_1 > \theta_0$. This reasoning gives us the following result (read carefully Definition 2.6.26 again).

Theorem 3.2.2. *Suppose the family of densities has monotone likelihood ratios. Then the test φ_t in (3.13) is uniformly most powerful for $H_0 : \theta \leq \theta_0$ versus $H_1 : \theta > \theta_0$ and has level $\alpha = \mathbb{E}_{\theta_0} \varphi_t$. Any $\alpha \in (0, 1)$ is possible.*

We discuss an example that is not an exponential family.

Example 3.2.3. *Suppose the data X_1, \dots, X_n are i.i.d. from the uniform distribution on $[0, \theta]$. The joint density $p_\theta(x)$ is positive if and only if $x_i \in [0, \theta]$ for $i = 1, \dots, n$ and this happens if and only if $M(x) = \min\{x_1, \dots, x_n\} \geq 0$ and $T(x) = \max\{x_1, \dots, x_n\} \leq \theta$. Thus*

$$p_\theta(x) = \begin{cases} 1/\theta^n & \text{if } M(x) \geq 0, T(x) \leq \theta \\ 0 & \text{otherwise.} \end{cases}$$

Suppose $\theta_2 > \theta_1$, $M(x) \geq 0$, and $T(x) \leq \theta_2$. Then

$$\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} = \begin{cases} (\theta_1/\theta_2)^n & \text{if } T(x) \leq \theta_1 \\ +\infty & \text{otherwise.} \end{cases}$$

This shows that the family of joint densities has monotone likelihood ratios. If we are interested in testing $\theta \leq 1$ versus $H_1 : \theta > 1$, the test function $\varphi_t(x)$ gives the UMP. This test has level

$$\mathbb{E}_1 \varphi_t(X) = \mathbb{P}_1(T \geq t) = 1 - t^n$$

and a specified level α can be achieved taking $t = (1 - \alpha)^{1/n}$. The power of this test is

$$\beta(\theta) = \mathbb{P}_\theta(T \geq t) = \begin{cases} 0 & \text{if } \theta < t \\ 1 - \frac{1-\alpha}{\theta^n} & \text{if } \theta \geq t. \end{cases}$$

3.2.2 Affine submodels*

Consider an exponential family as in (1.5). Consider a p -dimensional linear space $\mathcal{L} \subseteq \mathbb{R}^d$ and the corresponding $\Theta_0 = \Theta \cap \mathcal{L}$. There exists a simple affine change of coordinates from θ to (λ, ψ) such that $\Theta_0 = \{(\lambda, \psi) : \psi = \mathbf{0}\}$. If the basis of \mathcal{L} is given by the columns of $\mathbb{R}^{d \times p}$ then $\theta = A\lambda$ parametrizes Θ_0 . Thus, for $\theta \in \Theta_0$ we have

$$\langle \theta, t \rangle = \langle A\lambda, t \rangle = \langle \lambda, A^T t \rangle = \langle \lambda, u \rangle,$$

where $\mathbf{u} = A^\top \mathbf{t}$ is the sufficient statistics of the p -dimensional model parametrized by Θ_0 .

Thus, without loss of generality we assume $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ with $\boldsymbol{\psi} = \boldsymbol{\theta}_v$ and $\boldsymbol{\lambda} = \boldsymbol{\mu}_u$. We want to test the q -dimensional hypothesis ($q = d - p$) $\boldsymbol{\psi}_0 = \mathbf{0}$:

$$H_0 : \boldsymbol{\psi} = \mathbf{0}, \quad \text{versus} \quad \boldsymbol{\psi} \neq \mathbf{0}.$$

A simple example is that we have a regression model and want to delete a certain set of its regressors to reduce the dimension of the model. Alternatively, the primary model might be the smaller p -dimensional model with canonical statistics \mathbf{u} . We want to test if this model fits data, by embedding it in a wider model. As an example, consider the problem of testing a Gaussian graphical model with respect to graph G_0 (see Example 1.9.5) versus a Gaussian graphical model with respect to graph $G \supset G_0$.

We first consider a model reduction of the canonical statistics from $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ to \mathbf{u} . In Section 1.5 we argued that any inference on $\boldsymbol{\psi}$ should be done conditionally on \mathbf{u} . The argument was that \mathbf{u} provides no information about $\boldsymbol{\psi}$ as expressed in the likelihood factorisation (1.14) valid whether H_0 is true or not. All information provided by \mathbf{u} is consumed in estimating $\boldsymbol{\mu}_u$ and this parameter has no information about $\boldsymbol{\psi}$ by variational inference.

Recall the form of the distribution of $\mathbf{t} = (\mathbf{u}, \mathbf{v})$ in (1.6) and consider the conditional distribution for \mathbf{v} given \mathbf{u} . Inserting $\boldsymbol{\psi} = \mathbf{0}$ simplifies this to

$$f_0(\mathbf{v}|\mathbf{u}) = \frac{g(\mathbf{u}, \mathbf{v})}{g_0(\mathbf{u})}, \quad (3.14)$$

where index 0 indicates distribution under H_0 , and in particular $g_0(\mathbf{u}) = \int g(\mathbf{u}, \mathbf{v}) d\mathbf{v}$ is the structure function in the marginal exponential family for \mathbf{u} under H_0 . Note that, under the null, the conditional distribution of \mathbf{v} given \mathbf{u} is parameter free.

We propose the following test for $H_0 : \boldsymbol{\psi} = \mathbf{0}$ versus $H_1 : \boldsymbol{\psi} \neq \mathbf{0}$:

- (i) Use $f_0(\mathbf{v}|\bar{\mathbf{u}})$ as the test statistic, which, under the null, is equal to the conditional density of \mathbf{v} given $\bar{\mathbf{u}}$.
- (ii) Reject H_0 if $f_0(\bar{\mathbf{v}}|\bar{\mathbf{u}})$ is too small, and calculate the p -value as

$$\mathbb{P}\left(f_0(\mathbf{v}|\bar{\mathbf{u}}) \leq f_0(\bar{\mathbf{v}}|\bar{\mathbf{u}})\right) = \int_{f_0(\mathbf{v}|\bar{\mathbf{u}}) \leq f_0(\bar{\mathbf{v}}|\bar{\mathbf{u}})} f_0(\mathbf{v}|\bar{\mathbf{u}}) d\mathbf{v},$$

where $\bar{\mathbf{u}}, \bar{\mathbf{v}}$ denote the observed quantities.

These test follow the Fisher's principle of **exact tests**. "Exactness" comes from the fact that the test achieves exactly the desired size α (as opposed to approximate tests). One obvious question is about the power of such a procedure.

Consider a special case of the above setting when $\mathbf{t} = (\mathbf{u}, v)$ with v one-dimensional. In this case θ_v is one-dimensional too and we can obtain a uniform most powerful test for $H_0 : \psi = 0$ versus $H_1 : \psi \neq 0$. Recall that, under the null, the conditional distribution of v given \mathbf{u} is parameter free. We will then use directly v as the test statistic with \mathbf{u} fixed and consider the test $\varphi^*(x)$ defined as

$$\varphi^*(x) = \begin{cases} 1 & \text{if } v > c^*(\mathbf{u}), \\ 1 & \text{if } v < c_*(\mathbf{u}), \\ \gamma^*(\mathbf{u}) & \text{if } v = c^*(\mathbf{u}), \\ \gamma_*(\mathbf{u}) & \text{if } v = c_*(\mathbf{u}), \\ 0 & \text{if } v \in (c_*(\mathbf{u}), c^*(\mathbf{u})), \end{cases}$$

with $c(\cdot)$ and $\gamma(\cdot)$ adjusted so that the test has exactly level α .

Theorem 3.2.4. *If the exponential family is regular and θ_v is one dimensional, then φ^* is a uniformly most powerful unbiased test of $H_0 : \theta_v = 0$ versus $H_1 : \theta_v \neq 0$.*

Proof. See Theorem 13.6 in ¹. □

Example 3.2.5 (Mean value test for a normal distribution). *Consider testing $H_0 : \mu = 0$ for a sample of size n from $N(\mu, \sigma^2)$. In Example 1.1.4 we note that the canonical parameters of this family are $(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2})$ so equivalently we test $\theta_1 = 0$ with $v = \sum_i x_i$. The exact test should be conditional, given $u = \sum_i x_i^2 = \|x\|^2$. Thus we are free to eliminate the scale parameter by considering v/\sqrt{u} instead of v . Indeed, use Exercise 3.7.3 to conclude that, under the null, $u \perp\!\!\!\perp v/\sqrt{u}$ and the vector $\frac{1}{\sqrt{u}}(x_1, \dots, x_n)$ is uniformly distributed on the $(n-1)$ -dimensional unit sphere and so the distribution of v/\sqrt{u} is also parameter-free. Exercise 3.7.3 also allows to conclude that v/\sqrt{u} is independent of u . Thus we can forget about conditioning on u , that is, we may restrict attention to the marginal distribution of v/\sqrt{u} .*

The last step is to show that this test is equivalent to the t-test. Let $\tau = \sqrt{n}\bar{x}/s$ be the t-test statistic. Let us rewrite v/\sqrt{u} as

$$\frac{n\bar{x}}{\sqrt{\sum_i x_i^2}} = \frac{n\bar{x}}{\sqrt{(n-1)s^2 + n\bar{x}^2}} = \frac{n\bar{x}/s}{\sqrt{(n-1) + (\sqrt{n}\bar{x}/s)^2}} = \frac{\sqrt{n}\tau}{\sqrt{(n-1) + \tau^2}}.$$

The right-hand side is seen to be an odd function and monotone function of τ . Thus, a single tail or a symmetric pair of tails in u/\sqrt{v} is equivalent to a single or symmetric pair of tails in the usual t-test.

3.3 Sequential testing*

The aim of this section is to give a brief introduction into sequential testing. For simplicity we focus on the **sequential probability ratio**

¹ Robert W. Keener. *Theoretical statistics*. Springer Texts in Statistics. Springer, New York, 2010. Topics for a core course

← Exercise 3.7.3

← Exercise 3.7.4

test. This test was suggested by Wald for simple versus simple testing with *i.i.d* observations with optional stopping. Let X_1, X_2, \dots be *i.i.d* from a distribution with density $p_k, k = 0, 1$, and consider testing $H_0 : k = 0$ versus $H_1 : k = 1$. Define

$$\text{LR}_n = \text{LR}_n(X_1, \dots, X_n) = \prod_{i=1}^n \frac{p_1(X_i)}{p_0(X_i)}, \quad (3.15)$$

the likelihood ratio for the first n observations. By convention, $\text{LR}_0 = 1$. From the Neyman-Pearson theorem we know that for a fixed sample size the best test rejects H_0 according to the size of LR_n . In the **Sequential Probability Ratio Test** (SPRT) at each step the researcher has three options: stop and accept H_0 , stop and accept H_1 , or continue sampling. For the SPRT these options are resolved by comparing the likelihood ratio with two critical values $\gamma_0 < 1 < \gamma_1$ in the following manner:

1. If $\text{LR}_n \in (\gamma_0, \gamma_1)$, take another observation.
2. If $\text{LR}_n \geq \gamma_1$, reject H_0 .
3. If $\text{LR}_n \leq \gamma_0$, accept H_0 .

Formally, the sample size for this SPRT is then a random variable defined as

$$N := \inf\{n : \text{LR}_n \notin (\gamma_0, \gamma_1)\}.$$

We will set the thresholds to provide desired power

$$\beta = \mathbb{P}_1(\text{LR}_N \geq \gamma_1)$$

and size

$$\alpha = \mathbb{P}_0(\text{LR}_N \geq \gamma_1).$$

Note that both quantities involve N , which makes the analysis more subtle. To simplify the notation, for a fixed $n \in \mathbb{N}$, let $\mathbf{x} := (x_1, \dots, x_n)$ and write $p_k(\mathbf{x}) := \prod_{i=1}^n p_k(x_i), k = 0, 1$. Let $R_1 = \{\mathbf{x} : \text{LR}_n \geq \gamma_1\}$ be the rejection region. Note that

$$\mathbb{P}_1(\text{LR}_n \geq \gamma_1) = \int_{R_1} p_1(\mathbf{x}) d\mathbf{x} = \int_{R_1} \text{LR}_n p_0(\mathbf{x}) d\mathbf{x} \geq \gamma_1 \mathbb{P}_0(\text{LR}_n \geq \gamma_1). \quad (3.16)$$

Similarly, denoting $R_0 = \{\text{LR}_n \leq \gamma_0\}$,

$$\mathbb{P}_0(\text{LR}_n \leq \gamma_0) = \int_{R_0} p_0(\mathbf{x}) d\mathbf{x} = \int_{R_0} \text{LR}_n^{-1} p_1(\mathbf{x}) d\mathbf{x} \geq \gamma_0^{-1} \mathbb{P}_1(\text{LR}_n \leq \gamma_0). \quad (3.17)$$

With a bit more careful treatment (using Wald's likelihood ratio identity, Theorem 2.3.3 and Section 3.1.1.1 in ²) we can replace n with the

² Alexander Tartakovsky, Igor Nikiforov, and Michele Basseville. *Sequential analysis: Hypothesis testing and changepoint detection*. CRC Press, 2014

stopping time N both in (3.16) and (3.17). It follows that $\beta \geq \gamma_1 \alpha$ and $(1 - \alpha) \geq \gamma_0^{-1}(1 - \beta)$ and so

$$\gamma_1 \leq \frac{\beta}{\alpha} \quad \text{and} \quad \gamma_0 \geq \frac{1 - \beta}{1 - \alpha}. \quad (3.18)$$

The relation between (α, β) and (γ_0, γ_1) is complicated. The above inequalities lead to an approximate analysis. Suppose (α^*, β^*) are some desired size and the power. Using the proposal of Wald, we set $\gamma_1^* = \frac{\beta^*}{\alpha^*}$ and $\gamma_0^* = \frac{1 - \beta^*}{1 - \alpha^*}$. Note that α^*, β^* are typically not equal to the real size and power. The bounds in (3.18) guarantee however that if α^* is small and β^* is large then $\alpha \approx \alpha^*$ and $\beta \approx \beta^*$; see Figure 3.2.

Expected stopping time of SPRT. To gain insight into this issue, let us consider the expected stopping time. We can calculate the expected value of N as follows. First observe that, for any fixed time n ,

$$\mathbb{E}_k(\log \text{LR}_n) = \sum_{i=1}^n \mathbb{E}_k \left(\log \frac{p_1(X_i)}{p_0(X_i)} \right) = \begin{cases} nD(p_1 \| p_0) & \text{if } k = 1 \\ -nD(p_0 \| p_1) & \text{if } k = 0, \end{cases}$$

where $D(p_0 \| p_1)$ is the KL-divergence between p_0 and p_1 .

Proposition 3.3.1 (Wald's identity). *Let Y_1, Y_2, \dots be independent and identically distributed random variables with mean μ and suppose $\mathbb{E}|Y_i| < C$ for some C . Let N be any integer-valued random variable such that $E[N] < \infty$ and $\{N = n\} \in \sigma(Y_1, \dots, Y_n)$. Then $\mathbb{E}[\sum_{i=1}^N Y_i] = \mu \mathbb{E}[N]$.*

Proof. Start by noting that the event $\{N \geq i\} = (\cup_{j=1}^{i-1} \{N = j\})^c$. Thus, the event is independent of Y_i, Y_{i+1}, \dots (since it is determined by Y_1, \dots, Y_{i-1}). From this we see that

$$\sum_{i=1}^{\infty} \mathbb{E}(|Y_i| \mathbb{1}(N \geq i)) = \sum_{i=1}^{\infty} \mathbb{E}(|Y_i|) \mathbb{P}(N \geq i) \leq C \mathbb{E}N < \infty. \quad (3.19)$$

Write

$$\mathbb{E}\left(\sum_{i=1}^N Y_i\right) = \mathbb{E}\left(\sum_{i=1}^{\infty} \mathbb{1}(N \geq i) Y_i\right) = \sum_{i=1}^{\infty} \mathbb{E}(\mathbb{1}(N \geq i) Y_i).$$

(the interchange of expectation and summation is justified by the dominated convergence theorem and (3.19)). Therefore,

$$\sum_{i=1}^{\infty} \mathbb{E}[\mathbb{1}(N \geq i) Y_i] = \mathbb{E}(Y_1) \sum_{i=1}^{\infty} \mathbb{E}(\mathbb{1}(N \geq i)) = \mu \sum_{i=1}^{\infty} \mathbb{P}(N \geq i) = \mu \mathbb{E}(N).$$

□

So, by Wald's Identity we have

$$\mathbb{E}_k(\log \text{LR}_N) = \begin{cases} \mathbb{E}_1(N)D(p_1 \| p_0) & \text{if } k = 1 \\ -\mathbb{E}_0(N)D(p_0 \| p_1) & \text{if } k = 0. \end{cases} \quad (3.20)$$

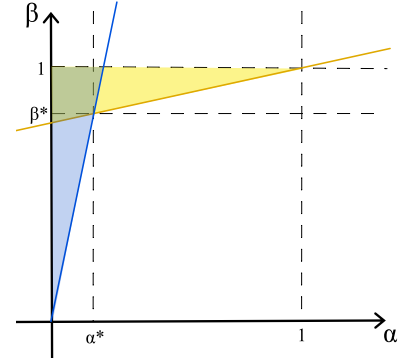


Figure 3.2: Bounds in (3.18) with $\gamma_0 = \frac{1 - \beta^*}{1 - \alpha^*}$ and $\gamma_1 = \frac{\beta^*}{\alpha^*}$.

Now to obtain an expression for $\mathbb{E}_k(N)$ we will derive another formula for $\mathbb{E}_k[\log \text{LR}_N]$. Let us assume the value of the likelihood ratio is approximately equal to a threshold level when the SPRT terminates. The value of the likelihood ratio will typically be just slightly greater/lower than the upper/lower threshold level. Using this approximation and smoothing (note that $\mathbb{P}_k(\text{LR}_N \in (\gamma_0, \gamma_1)) = 0$) we can write

$$\begin{aligned}\mathbb{E}_0(\log \text{LR}_N) &\approx \alpha \log(\gamma_1) + (1 - \alpha) \log(\gamma_0) \\ &\approx \alpha \log \frac{\beta}{\alpha} + (1 - \alpha) \log \frac{1 - \beta}{1 - \alpha},\end{aligned}$$

where we used the fact that $\gamma_1 = \frac{\beta}{\alpha}$ and $\gamma_0 = \frac{1 - \beta}{1 - \alpha}$. Similarly,

$$\mathbb{E}_1(\log \text{LR}_N) \approx \beta \log \frac{\beta}{\alpha} + (1 - \beta) \log \frac{1 - \beta}{1 - \alpha}. \quad (3.21)$$

Denoting by $\pi_\alpha \sim \text{Bern}(\alpha)$ and $\pi_\beta \sim \text{Bern}(\beta)$, we get $\mathbb{E}_0(\log \text{LR}_N) \approx -D(\pi_\alpha \| \pi_\beta)$ and $\mathbb{E}_1(\log \text{LR}_N) \approx D(\pi_\beta \| \pi_\alpha)$. With these approximations, using (3.20), we obtain expressions for $\mathbb{E}_k(N)$:

$$\begin{aligned}\mathbb{E}_0(N) &\approx \frac{D(\pi_\alpha \| \pi_\beta)}{D(p_0 \| p_1)} \\ \mathbb{E}_1(N) &\approx \frac{D(\pi_\beta \| \pi_\alpha)}{D(p_1 \| p_0)}.\end{aligned}$$

Note that the expected stopping times increase as the KL divergences decreases (as the two densities become less distinguishable). Increasing β or decreasing α also increases the expected stopping time.

Optimality of SPRT The expected stopping time of the SPRT that we determined above is optimal. No other test can achieve the same β and α with a smaller expected number of samples, under either hypothesis, as the following result shows.

Lemma 3.3.2 (Lower bound on expected stopping time of any testing procedure (Wald&Wolfowitz 1948)). *Let α and β be given and consider any sequential test with size $\leq \alpha$ and power $\geq \beta$. Then the expected stopping times N' for the test satisfy $\mathbb{E}_k N' \geq \mathbb{E}_k N$ for $k = 0, 1$, where N is the stopping time of the corresponding SPRT with size α and power β .*

The lemma shows that if no other test can have error levels as small or smaller than the SPRT and have expected stopping times less than the values computed above for the SPRT.

Example 3.3.3 (Sequential testing in Gaussian case). *Let X_1, X_2, \dots be an i.i.d. sequence of normal variables $N(\mu, 1)$. Consider the simple binary testing problem: $H_0 : \mu = 0, H_1 : \mu = \mu_0 > 0$. For simplicity, let us specify equal probabilities of error, that is, $\alpha = 1 - \beta \ll \frac{1}{2}$. In this case the*

optimal cut-off for the (non-sequential) likelihood ratio test is to reject the null if $\log \text{LR}_n \geq 0$. It also easily follows that

$$\alpha = \Phi \left(\frac{\sqrt{n}\mu_0}{2} \right),$$

where Φ is the CDF of the standard normal variable. So the number of samples required for a specified α is

$$n = \frac{2(\Phi^{-1}(\alpha))^2}{\mu_0}.$$

Since $D(p_0||p_1) = \mu^2/2$ and $D(\pi_\alpha||\pi_\beta) = D(\pi_\beta||\pi_\alpha) = (1 - 2\alpha) \log \frac{1-\alpha}{\alpha}$, the expected stopping time of the SPRT in this case is approximately

$$\mathbb{E}_0(N) = \mathbb{E}_1(N) = \frac{2(1-2\alpha)}{\mu_0^2} \log \frac{1-\alpha}{\alpha}.$$

Compare the two quantities to see that the sequential sample requirement is indeed preferred.

3.4 Motivating multiple testing

Selective inference means searching for interesting patterns in data, with statistical guarantees that account for the search process. It encompasses multiple testing, post-selection inference, and adaptive or interactive inference. There are two main situations for considering simultaneous inference. One situation is when we test multiple hypothesis. Another is when we try to obtain simultaneous coverage for confidence intervals for a multiple parameters. In this section we will focus on multiple testing.

As a general motivation, consider the general problem of simultaneously testing a finite number of null hypotheses H_0^i for $i = 1, \dots, m$.

Example 3.4.1. *Suppose that we have m genes and data about expression levels for each gene among healthy individuals and those with lung cancer.*

	Healthy (k patients)	Lung cancer (l patients)
Expression Level of Gene i	$x_{ij}^{(0)}, 1 \leq j \leq k$	$x_{ij}^{(1)}, 1 \leq j \leq l$

The i -th null hypothesis, denoted H_0^i , would state that the mean expression level of the i -th gene is the same in both groups of patients.

We assume that the tests for the individual hypotheses are available (T_i test statistic, R_i rejection region, P_i the associated p-value, i.e. the smallest α leading to rejection) and the problem is how to combine them into a simultaneous test procedure. The easiest yet extremely naive approach (as illustrated in the webcomic xkcd³) is to disregard the multiplicity and simply test each hypothesis at level α . However, with such a procedure the probability of one or more false rejections rapidly increases with n . For example, if all test are independent of size α and all null hypotheses are true then

³ See <https://xkcd.com/882/>

$$\mathbb{P}_0(\bigcap_{i=1}^m \{T_i \notin R_i\}) = \prod_{i=1}^m \mathbb{P}_0(T_i \notin R_i) = (1 - \alpha)^m.$$

In this sense the claim that the procedure controls the probability of false rejections at level α is clearly misleading. A similar situation emerges when constructing a confidence region for a parameter vector using individual confidence intervals for each component.

3.5 Family-wise error rate

Let $\mathcal{H}_0 \subseteq \{1, \dots, m\}$ be the index set of the true hypotheses and let $\mathcal{R} \subseteq \{1, \dots, m\}$ be the set of rejected hypotheses. Denote $m_0 = |\mathcal{H}_0|$. The **family-wise error rate (FWER)** is

$$\text{FWER} = \mathbb{P}(|\mathcal{H}_0 \cap \mathcal{R}| \geq 1). \tag{3.22}$$

A natural approach is to replace the usual condition for testing a single hypothesis, that the probability of false rejection not exceed α , by the requirement

$$\text{FWER} \leq \alpha$$

for all possible combinations of true and false hypotheses. Methods that control the FWER are often described by the p-values of individual tests.

3.5.1 Example: Gaussian sequence model

For simplicity of the discussion, for most of this section, we restrict our discussion to the important example given by the **Gaussian sequence model**. Consider a model $Y_i = \mu_i + \varepsilon_i$ for $i = 1, \dots, m$, where $\varepsilon_i \sim N(0, 1)$. For now we will not assume that ε_i are independent⁴. Also the case, when variance of the noise is a general $\sigma^2 > 0$ but known can be easily covered. The typical question that is asked about this model is how we can test, which elements of $\mu = (\mu_1, \dots, \mu_m)$ are zero, or perhaps, which elements of μ are equal to each other. There are many other questions that can be phrased as linear equalities or inequalities in the vector μ .

⁴With independence, this model already appeared in our discussion of the Stein's paradox.

In the question of testing which μ_i are zero, we already mentioned the naive approach: take $z_{\alpha/2} = \Phi^{-1}(1 - \alpha/2)$. Test $H_0^i : \mu_i = 0$ using the test statistic $T_i = |Y_i|$ and rejection rule $\mathbb{1}\{|Y_i| > z_{\alpha/2}\}$. One classical fix is given by the Bonferroni correction, which is to use α/m instead of α . Now, we test $H_0^i : \mu_i = 0$ with the test $\mathbb{1}\{|Y_i| > z_{\alpha/2m}\}$. Assuming all the nulls are true, we obtain

$$\text{FWER} = \mathbb{P}_0(\exists i \mid \varepsilon_i| > z_{\alpha/2m}) \leq \sum_{i=1}^m \mathbb{P}_0(|\varepsilon_i| > z_{\alpha/2m}) = m \frac{\alpha}{m} = \alpha,$$

where in the inequality we used the union bound. It is clear that the union bound can be conservative, which implies that the Bonferroni bound can be conservative too. In the special case when all ε_i are independent, the events $\{|\varepsilon_i| > z_{\alpha/2m}\}$ are independent and

$$\begin{aligned} \text{FWER} &= \mathbb{P}_0(\exists i \mid \varepsilon_i| > z_{\alpha/2m}) = 1 - \prod_{i=1}^m \mathbb{P}_0(|\varepsilon_i| \leq z_{\alpha/2m}) \\ &= 1 - \left(1 - \frac{\alpha}{m}\right)^m \approx 1 - e^{-\alpha} \approx \alpha. \end{aligned}$$

Therefore, in this case, the Bonferroni procedure provides a good control over the family-wise error rate. This also tells us that if we have many hypotheses (e.g. $m = 10,000$ genes in the biological example) then Bonferroni's test has size approximately $1 - e^{-\alpha}$, which for small α is approximately α . For example, if $\alpha = 0.05$, then $1 - e^{-\alpha} = 0.04877 \dots$. So to get a test of size 0.05, we could test each hypothesis at level $0.0512/m$.

Remark 3.5.1 (Šidák correction). *The above calculation shows that we could use $z_{\tilde{\alpha}/2}$ with $\tilde{\alpha} = 1 - (1 - \alpha)^{1/m}$ in order to get the error rate precisely α under independence.*

The Bonferroni correction can become overly conservative in the case when ε_i are dependent. In the extreme situation, when they are all equal we get

$$\text{FWER} = \mathbb{P}(\exists i \ |\varepsilon_i| > z_{\alpha/2m}) = \mathbb{P}(|\varepsilon_1| > z_{\alpha/2m}) = \frac{\alpha}{m}.$$

As we said, one of the problems of selective inference can be to find guarantees for simultaneous coverage. In what follows we are going to assume full independence of the errors ε_i . Consider a problem in which for all $\mathbf{v} \in \mathbb{R}^m$ we want to build a confidence interval $\text{Cl}_{\mathbf{v}}$ to cover the parameter $\mathbf{v}^\top \mu$ with the property that

$$\mathbb{P}(\mathbf{v}^\top \mu \in \text{Cl}_{\mathbf{v}} \text{ for all } \mathbf{v} \in \mathbb{R}^m) \geq 1 - \alpha.$$

Since there are infinitely many vectors \mathbf{v} , using the Bonferroni correction is not possible. One way to solve this problem is to use the **Scheffé's method**. First, it is clear that the problem depends only on the direction of \mathbf{v} and not on its norm. Assume then that $\|\mathbf{v}\| = 1$. The confidence interval $\text{Cl}_{\mathbf{v}}$ will be centered around $\mathbf{v}^\top Y$. To get its length, we need to bound $|\mathbf{v}^\top Y - \mathbf{v}^\top \mu|$ for all \mathbf{v} such that $\|\mathbf{v}\| = 1$. We use

$$\chi_m(\alpha) = \text{the } (1 - \alpha)\text{-quantile of the } \chi_m \text{ distribution}$$

and define

$$\text{Cl}_{\mathbf{v}} = (\mathbf{v}^\top Y - \chi_m(\alpha), \mathbf{v}^\top Y + \chi_m(\alpha)).$$

We have

$$\begin{aligned} \mathbb{P}(\mathbf{v}^\top \mu \in \text{Cl}_{\mathbf{v}} \text{ for all } \|\mathbf{v}\| = 1) &= \mathbb{P}(|\mathbf{v}^\top Y - \mathbf{v}^\top \mu| < \chi_m(\alpha) \text{ for all } \|\mathbf{v}\| = 1) \\ &= \mathbb{P}(\|Y - \mu\| < \chi_m(\alpha)) = \mathbb{P}(\|\varepsilon\| < \chi_m(\alpha)) \\ &= \alpha. \end{aligned}$$

Remark 3.5.2. *In this section, we always assumed that $\sigma^2 = 1$ or equivalently that σ^2 is known. If it is unknown, it can be still often estimated using replicates. So suppose $Y_{ij} = \mu_i + \varepsilon_{ij}$, where $\varepsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ for $i = 1, \dots, m$ and $j = 1, \dots, n$. Let $\bar{Y}_i = \frac{1}{n} \sum_j Y_{ij}$, $\bar{\varepsilon}_i = \frac{1}{n} \sum_j \varepsilon_{ij}$. Then $\bar{Y}_i = \mu_i + \bar{\varepsilon}_i$ with $\bar{\varepsilon}_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2/n)$. Using ideas in Example 1.4.5, we can then simply construct an estimate $\hat{\sigma}^2$ of σ^2 , where $\hat{\sigma}^2 \perp\!\!\!\perp \bar{\varepsilon}_1, \dots, \bar{\varepsilon}_m$. Moreover,*

$$\frac{\hat{\sigma}^2}{\sigma^2} (n-1)m \sim \chi_{(n-1)m}^2.$$

Now the correct modification of Bonferroni bounds is to use the quantiles of $t_{(n-1)m}$ instead of $N(0, 1)$. In the Scheffé's method we use $F_{m, (n-1)m}$ in place of χ_m .

3.5.2 Bonferroni and Holm

First, recall a basic fact of statistical tests. Suppose a null hypothesis H_0 is true, and we perform a statistical test of H_0 and obtain a p-value P . What is the distribution of P ? Recall that the p-value is the probability of obtaining test results at least as extreme as the result actually observed, under the assumption that the null hypothesis is correct. If our test statistic T has a continuous distribution under H_0 with CDF F , and the rejection regions are of the form $\{X \in \mathcal{X} : T(X) \leq t_\alpha\}$, then the p-value is just the lower tail probability $P = F(T)$. Thus, for any $u \in (0, 1)$

$$\mathbb{P}_0(P \leq u) = \mathbb{P}_0(F(T) \leq u) = \mathbb{P}_0(T \leq F^{-1}(u)) = F(F^{-1}(u)) = u.$$

So $P \sim U(0, 1)$ under the null. Similarly, $P \sim U(0, 1)$ if we reject for large T , or both large and small T .

If T has a discrete distribution under H_0 , then so does P , so the null distribution of P would not be exactly uniform. However, we still have that

$$\mathbb{P}(P \leq u) \leq u \quad \text{for all } u \in (0, 1). \quad (3.23)$$

To show this define

$$G(u) = \sup\{y : F(y) \leq u\}. \quad (3.24)$$

Note that, by definition, $F(t) \leq u$ implies that $t \leq G(u)$. Since F is non-decreasing and right-continuous we have $F(G(u)) \geq u$ ⁵. Suppose that u is such that $F(G(u)) = u$, then

$$\mathbb{P}_0(P \leq u) = \mathbb{P}_0(F(T) \leq u) \leq \mathbb{P}(T \leq G(u)) = F(G(u)) = u.$$

If $F(G(u)) > u$, F had a jump at $\lambda = G(u)$

$$\lim_{y \rightarrow \lambda^-} F(y) \leq u, \quad \lim_{y \rightarrow \lambda^+} F(y) > u.$$

By Remark C.1.3 the function $\mathbb{P}(T < \lambda)$ is left-continuous in λ . Thus, we have

$$\mathbb{P}_0(P \leq u) = \mathbb{P}_0(F(T) \leq u) = \mathbb{P}(T < \lambda) = \lim_{t \rightarrow \lambda^-} \mathbb{P}(T < t) \leq \lim_{y \rightarrow \lambda^-} F(y) \leq u.$$

This gives us one simple way of dealing with multiple testing problem.

Theorem 3.5.3 (Bonferroni Procedure). *If, for $i = 1, \dots, m$, hypothesis H_0^i is rejected when $P_i \leq \alpha/m$, then the FWER for the simultaneous testing of H_0^1, \dots, H_0^m satisfies $\text{FWER} \leq \alpha$.*

⁵ Indeed, if $t > G(u)$ then $F(t) > u$. By right-continuity, if converging to $G(u)$ from the right gives $F(G(u)) \geq u$.

Proof. Suppose hypotheses H_0^i with $i \in \mathcal{H}_0$ are true and the other are false. From the union bound it follows that

$$\begin{aligned} \text{FWER} &= \mathbb{P}(\text{reject any } H_0^i \text{ with } i \in \mathcal{H}_0) \leq \sum_{i \in \mathcal{H}_0} \mathbb{P}(\text{reject } H_0^i) \\ &= \sum_{i \in \mathcal{H}_0} \mathbb{P}(P_i \leq \frac{\alpha}{m}) \stackrel{(3.23)}{\leq} \sum_{i \in \mathcal{H}_0} \frac{\alpha}{m} \leq \frac{m_0}{m} \alpha \leq \alpha. \end{aligned}$$

□

From the proof, it is clear that, although this procedure controls FWER, it is generally too conservative to be useful in detecting H_0^i that are false. If m_0 is smaller than m , using α/m_0 would be preferable. The problem, of course, is that m_0 is not known. Note however that if one H_0^i is false then $m_0 \leq m - 1$ and we could use $\alpha/(m - 1)$ as the threshold. If this H_0^i was true and we rejected it then we already made a mistake so we can do whatever we want and there is no harm in using $\alpha/(m - 1)$ for the other hypotheses (if we make a mistake, it does not matter how many).

The Holm procedure tries to make use of the above observations. It can conveniently be stated in terms of the p-values P_1, \dots, P_n of the n individual tests. The procedure starts by sorting the p-values. Call them $P_{(1)} \leq \dots \leq P_{(m)}$. Then it goes as follows:

1. If $P_{(1)} \leq \frac{\alpha}{m}$, reject $H_0^{(1)}$ and continue. Else stop.
2. If $P_{(2)} \leq \frac{\alpha}{m-1}$, reject $H_0^{(2)}$ and continue. Else stop.
- ...
- m. If $P_{(m)} \leq \alpha$, reject $H_0^{(m)}$.

In other words, the procedure finds the smallest r such that $P_{(r)} > \frac{1}{m-(r-1)}\alpha$ and it rejects the null hypotheses $H_0^{(1)}, \dots, H_0^{(r-1)}$.

Theorem 3.5.4. *The Holm procedure satisfies $\text{FWER} \leq \alpha$.*

Proof. Suppose \mathcal{H}_0 is the set of true hypotheses. Order the P-values as above $P_{(1)} \leq \dots \leq P_{(m)}$. Ideally, if $i \in \mathcal{H}_0$ then P_i appears later in this order sequence. Let j be the smallest (random) index satisfying

$$P_{(j)} = \min_{i \in \mathcal{H}_0} P_i.$$

Note that, by construction, $j \leq m - m_0 + 1$ (there are at least $m_0 - 1$ indices following (j)). Now, the Holm procedure commits a false rejection if $r \geq j + 1$, or in other words,

$$P_{(1)} \leq \frac{\alpha}{m}, \quad P_{(2)} \leq \frac{\alpha}{m-1}, \quad \dots, \quad P_{(j)} \leq \frac{\alpha}{m-j+1}$$

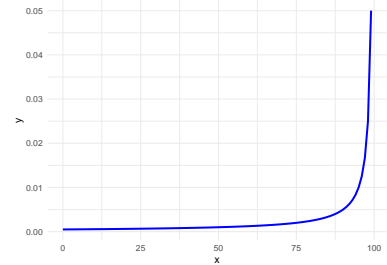


Figure 3.3: The plot of the function $\frac{\alpha}{m-x}$ with $m = 100$ and α . The ordered p-values $P_{(1)}, \dots, P_{(r)}$ will lie below this curve.

which implies that

$$\min_{i \in \mathcal{H}_0} P_i = P_{(j)} \leq \frac{\alpha}{m-j+1} \leq \frac{\alpha}{m_0}.$$

(This may be confusing because you may think that we could make a mistake by rejecting some other null. Note however that, if $P_{(j)} > \frac{\alpha}{m-j+1}$ then the procedure stops and no other mistakes are done!) We thus get

$$\begin{aligned} \mathbb{P}(|\mathcal{H}_0 \cap \mathcal{R}| \geq 1) &= \mathbb{P}(P_{(1)} \leq \frac{\alpha}{m}, \dots, P_{(j)} \leq \frac{\alpha}{m-j+1}) \\ &\leq \mathbb{P}(\min_{i \in \mathcal{H}_0} P_i \leq \frac{\alpha}{m-j+1}) \leq \mathbb{P}(\min_{i \in \mathcal{H}_0} P_i \leq \frac{\alpha}{m_0}). \end{aligned}$$

By the union bound, the probability of a false rejection is bounded above by

$$\mathbb{P}\left(\min_{i \in \mathcal{H}_0} P_i \leq \frac{\alpha}{m_0}\right) \leq \sum_{i \in \mathcal{H}_0} \mathbb{P}(P_i \leq \frac{\alpha}{m_0}) \leq \alpha.$$

□

This means that the Holm procedure is strictly more powerful than Bonferroni without any extra assumptions. However, Figure 3.3 also suggest that, when m is very large, the Holm procedure may not substantially differ from the Bonferroni correction.

3.6 False discovery rate

Controlling the FWER may be too conservative and greatly reduce our power to detect real effects, especially when m is large. In many modern “large-scale testing” applications, focus has shifted from FWER in (3.22) to the false-discovery proportion (FDP)

$$\text{FDP} = \frac{|\mathcal{H}_0 \cap \mathcal{R}|}{|\mathcal{R}| \vee 1}$$

and on procedures that control its expected value $\mathbb{E}(\text{FDP})$, called the **false-discovery rate (FDR)**. The FDR can be interpreted as follows: if $\text{FDR} \leq 0.1$, we expect around 90% of the discoveries to be true.

Controlling FDR is a shift in paradigm - we are willing to tolerate some type I errors (false discoveries), as long as most of the discoveries we make are still true. It has been argued that in applications where the statistical test is thought of as providing a “definitive answer” for whether an effect is real, FWER control is still the correct objective. In contrast, for applications where the statistical test identifies candidate effects that are likely to be real and which merit further study, it may be better to target FDR control.

3.6.1 Benjamini-Hochberg procedure

It is clear that the false discovery rate depends heavily on the number of true hypotheses. Thus, any procedure that controls FDR should be adaptive. As we showed earlier, if $i \in \mathcal{H}_0$ then $\mathbb{P}(P_i \leq u) \leq u$ for all $u \in (0, 1)$. The Benjamini-Hochberg (BH) procedure compares the sorted p-values to a diagonal cutoff line, finds the largest p-value that still falls below this line, and rejects the null hypotheses for the p-values up to and including this one. Formally, the BH procedure at level α is defined as follows:

1. Sort the p-values. Call them $P_{(1)} \leq \dots \leq P_{(m)}$ as before.
2. Find the largest r such that $P_{(r)} \leq \frac{r}{m}\alpha$.
3. Reject the null hypotheses $H_0^{(1)}, \dots, H_0^{(r)}$.

Remark 3.6.1. Just to avoid confusion with the Holm procedure, note that we do not require that $P_{(j)} \leq \frac{j}{m}\alpha$ for $j < r$.

It is useful to observe the following.

Lemma 3.6.2. Suppose the B-H procedure rejects exactly r hypotheses. Then $i \in \mathcal{R}$ if and only if $P_i \leq \frac{r}{m}\alpha$.

Proof. If $i \in \mathcal{R}$ then $P_i \leq P_{(r)} \leq \frac{r}{m}\alpha$, which proves the right implication. For the left implication we argue using a contrapositive statement. Suppose $i \notin \mathcal{R}$, let s be such that $P_i = P_{(s)}$. Then $s > r$ and $P_{(s)} > \frac{s}{m}\alpha > \frac{r}{m}\alpha$. \square

Theorem 3.6.3 (Benjamini and Hochberg). Consider tests of m null hypotheses. If the test statistics (or equivalently, p-values) of these tests are independent, then the FDR of the above procedure satisfies

$$\text{FDR} \leq \alpha \frac{m_0}{m} \leq \alpha.$$

Proof. We have

$$\text{FDR} = \mathbb{E} \left(\frac{|\mathcal{H}_0 \cap \mathcal{R}|}{|\mathcal{R}| \vee 1} \right) = \mathbb{E} \left(\frac{\sum_{i \in \mathcal{H}_0} \mathbb{1}\{i \in \mathcal{R}\}}{|\mathcal{R}| \vee 1} \right) = \sum_{i \in \mathcal{H}_0} \mathbb{E} \left(\frac{\mathbb{1}\{i \in \mathcal{R}\}}{|\mathcal{R}| \vee 1} \right).$$

Let $R = |\mathcal{R}|$ and let $R_{\setminus i}$ be the number of rejections if we replace P_i with zero and run BH_α on $P_1, \dots, P_{i-1}, 0, P_{i+1}, \dots, P_m$. Note that $P_i \perp\!\!\!\perp R_{\setminus i}$ and $R_{\setminus i} \geq 1$. Consider the following:

$$\text{Claim: } i \in \mathcal{R} \Leftrightarrow R = R_{\setminus i} \Leftrightarrow P_i \leq \frac{\alpha R_{\setminus i}}{m}.$$

Before we prove the claim, we show how it helps us to prove the

theorem. Assuming that the claim holds

$$\begin{aligned} \text{FDR} &= \sum_{i \in \mathcal{H}_0} \mathbb{E} \left(\frac{\mathbb{1}\{P_i \leq \frac{\alpha R_{\setminus i}}{m}\}}{R_{\setminus i}} \right) = \sum_{i \in \mathcal{H}_0} \mathbb{E} \left(\frac{\mathbb{P}(P_i \leq \frac{\alpha R_{\setminus i}}{m} | R_{\setminus i})}{R_{\setminus i}} \right) \\ &\leq \sum_{i \in \mathcal{H}_0} \mathbb{E} \left(\frac{\frac{\alpha R_{\setminus i}}{m}}{R_{\setminus i}} \right) = \alpha \frac{m_0}{m} \leq \alpha. \end{aligned}$$

It remains to prove the claim. We do it in several steps.

$[i \in \mathcal{R} \Rightarrow P_i \leq \frac{\alpha R_{\setminus i}}{m}]$: First note that BH_α is monotone in the p-values: If $P_i \geq P'_i$ for all then $R \leq R'$. From this it follows that $R \leq R_{\setminus i}$. Suppose $i \in \mathcal{R}$. Then

$$P_i \leq \frac{\alpha R}{m} \leq \frac{\alpha R_{\setminus i}}{m}.$$

$[P_i \leq \frac{\alpha R_{\setminus i}}{m} \Rightarrow R = R_{\setminus i}]$: By definition of $R_{\setminus i}$, we must have $R_{\setminus i}$ many values from $P_1, \dots, P_{i-1}, 0, P_{i+1}, \dots, P_m$ which are $\leq \frac{\alpha R_{\setminus i}}{m}$. Since $P_i \leq \frac{\alpha R_{\setminus i}}{m}$, it follows that there are $R_{\setminus i}$ many values from $P_1, \dots, P_{i-1}, P_i, P_{i+1}, \dots, P_m$ which are $\leq \frac{\alpha R_{\setminus i}}{m}$. In particular, $R \geq R_{\setminus i}$. The other inequality was shown above and so we get equality.

$[R = R_{\setminus i} \Rightarrow i \in \mathcal{R}]$ We prove this implication by contradiction. Suppose $R = R_{\setminus i}$ but $i \notin \mathcal{R}$. If $i \notin \mathcal{R}$ then $P_i > \frac{\alpha R}{m}$ and there are R many P_j 's with $P_j \leq \frac{\alpha R}{m}$. Now run BH_α on $P_1, \dots, P_{i-1}, 0, P_{i+1}, \dots, P_m$. Then there are $R + 1$ many values $\leq \frac{\alpha R}{m} \leq \frac{\alpha(R+1)}{m}$. Thus, we must have $R_{\setminus i} \geq R + 1 > R$, which gives contradiction. □

← Exercise 3.7.5

3.6.2 Benjamini-Yekutieli justification

The main problem with the validation of the Benjamini-Hochberg procedure is that it requires that the hypotheses are independent. Handling uniformly all dependence structure between the hypotheses is hard. One popular setting where the B-H procedure can be still validated is when the dependence between the P-values satisfies a form of positive dependence. We will now discuss this in more detail.

A set $S \subset \mathbb{R}^n$ is nondecreasing if $\mathbf{x} \in S$ and $\mathbf{y} \geq \mathbf{x}$ implies that $\mathbf{y} \in S$.

Definition 3.6.4. A random vector $\mathbf{X} = (X_1, \dots, X_m)$ is positively regression dependent on $I \subseteq \{1, \dots, m\}$ (PRDS) if $\mathbb{P}(\mathbf{X} \in S | X_i = x_i)$ is non-decreasing in x_i for every nondecreasing set S and any $i \in I$.

If (X_1, \dots, X_m) is PDRS on I and $Y_i := f_i(X_i)$ for all $1 \leq i \leq m$ with f_i strictly increasing or decreasing, then (Y_1, \dots, Y_m) is PRDS on I as well. Transformation of this form are called co-monotone transformations. Thus PRDS property is preserved under co-monotone

transformations. It follows that $P_i = \mathbb{P}(T_i \leq X_i)$ is PDRS as well (at least in the continuous case).

Theorem 3.6.5. *If the joint distribution of $\mathbf{P} = (P_1, \dots, P_m)$ is PRDS on the subset \mathcal{H}_0 , the Benjamini-Hochberg procedure controls the FDR at level $\frac{m_0}{m} \alpha$.*

The proof of this result relies on the following proposition.

Proposition 3.6.6. *If \mathbf{P} is PRDS on the set of true nulls, then the function $\mathbb{P}(\mathbf{P} \in S | P_i \leq t)$ for $i \in \mathcal{H}_0$ is non-decreasing in t for S a non-decreasing set.*

Proof. For any t , $\mathbb{P}(\mathbf{P} \in S | P_i \leq t) = \frac{\mathbb{P}(\mathbf{P} \in S, P_i \leq t)}{\mathbb{P}(P_i \leq t)}$. For $t > t'$, we have that

$$\mathbb{P}(\mathbf{P} \in S | P_i \leq t') = \frac{\mathbb{P}(\mathbf{P} \in S, P_i \leq t) + \mathbb{P}(\mathbf{P} \in S, P_i \in (t, t'])}{\mathbb{P}(P_i \leq t) + \mathbb{P}(P_i \in (t, t'])}$$

To show that $\mathbb{P}(\mathbf{P} \in S | P_i \leq t') \leq \mathbb{P}(\mathbf{P} \in S | P_i \leq t)$, it suffices to show that

$$\frac{\mathbb{P}(\mathbf{P} \in S, P_i \leq t)}{\mathbb{P}(P_i \leq t)} \leq \frac{\mathbb{P}(\mathbf{P} \in S, P_i \in (t, t'])}{\mathbb{P}(P_i \in (t, t'])}.$$

The last statement is because for any positive number a, b, c, d , $\frac{a}{b} \leq \frac{a+c}{b+d}$ if and only if $\frac{a}{b} \leq \frac{c}{d}$. If F_i denotes the CDF of P_i then

$$\begin{aligned} \mathbb{P}(\mathbf{P} \in S, P_i \in (t, t']) &= \mathbb{E} \mathbb{1}\{\mathbf{P} \in S, P_i \in (t, t']\} \\ &= \mathbb{E}[\mathbb{1}\{P_i \in (t, t']\} \mathbb{E}(\mathbb{1}\{\mathbf{P} \in S\} | P_i)] \\ &= \mathbb{E}[\mathbb{1}\{P_i \in (t, t']\} \mathbb{P}(\mathbf{P} \in S | P_i)] \\ &= \int_t^{t'} \mathbb{P}(\mathbf{P} \in S | P_i = s) dF_i(s) \\ &\stackrel{(PRDS)}{\geq} \int_t^{t'} \mathbb{P}(\mathbf{P} \in S | P_i = t) dF_i(s) \\ &= \mathbb{P}(\mathbf{P} \in S | P_i = t) \mathbb{P}(P_i \in (t, t']). \end{aligned}$$

Similarly we have

$$\begin{aligned} \mathbb{P}(\mathbf{P} \in S, P_i \leq t) &= \int_0^t \mathbb{P}(\mathbf{P} \in S | P_i = s) dF_i(s) \\ &\stackrel{(PRDS)}{\leq} \int_0^t \mathbb{P}(\mathbf{P} \in S | P_i = t) dF_i(s) \\ &= \mathbb{P}(\mathbf{P} \in S | P_i = t) \mathbb{P}(P_i \leq t). \end{aligned}$$

We have just shown that

$$\frac{\mathbb{P}(\mathbf{P} \in S, P_i \leq t)}{\mathbb{P}(P_i \leq t)} \leq \mathbb{P}(\mathbf{P} \in S | P_i = t) \leq \frac{\mathbb{P}(\mathbf{P} \in S, P_i \in (t, t'])}{\mathbb{P}(P_i \in (t, t'])}$$

as claimed. \square

Proof of Theorem 3.6.5. In the proof of Theorem 3.6.3 we noted that

$$\text{FDR} = \sum_{i \in \mathcal{H}_0} \mathbb{E} \left(\frac{\mathbb{1}\{i \in \mathcal{R}\}}{|\mathcal{R}| \vee 1} \right)$$

Note that it is enough to show that $\mathbb{E} \left(\frac{\mathbb{1}\{i \in \mathcal{R}\}}{|\mathcal{R}| \vee 1} \right) \leq \alpha/m$ for $i \in \mathcal{H}_0$.

By Lemma 3.6.2, if r rejections are made then H_0^i is rejected if and only if $P_i \leq \frac{\alpha r}{m}$. Hence, $\mathbb{1}\{i \in \mathcal{R}\} = \mathbb{1}\{P_i \leq \frac{\alpha r}{m}\}$. This gives

$$\mathbb{E} \left(\frac{\mathbb{1}\{i \in \mathcal{R}\}}{R \vee 1} \right) = \mathbb{E} \left(\mathbb{E} \left[\frac{\mathbb{1}\{i \in \mathcal{R}\}}{R \vee 1} \mid R \right] \right) = \mathbb{E} \left(\frac{1}{R \vee 1} \mathbb{P}[i \in \mathcal{R} \mid R] \right) = \sum_{r=1}^m \frac{1}{r} \mathbb{P}(P_i \leq \frac{\alpha r}{m}, R = r).$$

For any true null we now have

$$\begin{aligned} \sum_{r=1}^m \frac{1}{r} \mathbb{P}(P_i \leq \frac{\alpha r}{m}, R = r) &= \sum_{r=1}^m \frac{1}{r} \mathbb{P}(P_i \leq \frac{\alpha r}{m}) \mathbb{P}(|\mathcal{R}| = r \mid P_i \leq \frac{\alpha r}{m}) \\ &\leq \frac{\alpha}{m} \sum_{r=1}^m \mathbb{P}(R = r \mid P_i \leq \frac{\alpha r}{m}). \end{aligned}$$

Hence, it suffices to show that $\sum_{r=1}^m \mathbb{P}(R = r \mid P_i \leq \frac{\alpha r}{m}) \leq 1$. Observe that $\{R \leq r\}$ is an increasing event, that is, it can be written as $\{\mathbf{P} \in S\}$ for some non-decreasing set S . This is because increasing all p-values increases the p-value at each rank. Hence, any ranked p-value above the threshold remains above its threshold, that is, we accept at least as many as before and hence, do not reject more hypotheses. Using this, we get that

$$\begin{aligned} \sum_{r=1}^m \mathbb{P}(R = r \mid P_i \leq \frac{\alpha r}{m}) &= (\mathbb{P}(R \leq m \mid P_i \leq \alpha) - \mathbb{P}(R \leq 0 \mid P_i \leq \frac{\alpha}{m})) \\ &\quad + \sum_{r=1}^{m-1} \left(\mathbb{P}(R \leq r \mid P_i \leq \frac{\alpha r}{m}) - \mathbb{P}(R \leq r \mid P_i \leq \frac{\alpha(r+1)}{m}) \right). \end{aligned}$$

Note that each summand in the second line is non-positive since $\mathbb{P}(R \leq r \mid P_i \leq x)$ is increasing in x . Also $\mathbb{P}(R \leq 0 \mid P_i \leq \frac{\alpha}{m}) \geq 0$, which implies that

$$\sum_{k=1}^m \mathbb{P}(R = k \mid P_i \leq \frac{\alpha k}{m}) \leq \mathbb{P}(R \leq m \mid P_i \leq \alpha) \leq 1.$$

As stated before, this proves the upper bound on the FDR. \square

3.7 Exercises

Exercise 3.7.1. Show that if an α -level test maximizes $\mathbb{E}_1 \varphi$ then it must be essentially equal to the likelihood ratio test. (Hint: Use the proof of Theorem 3.1.1)

Exercise 3.7.2. Show that the threshold $\lambda = 0$ in Example 3.1.4 is optimal to optimize the Bayes risk with equal weights on each class. Hint: Note that $Z = \Sigma^{-1}(X - \mu_0)$ is standard normal if X comes from the first class.

Exercise 3.7.3. Show that if Z is d -dimensional standard Gaussian vector then $D = \|Z\|$ and $U = Z/\|Z\|$ are independent. Conclude that every d -dimensional vector $X \sim N_d(\mu, \Sigma)$ admits a stochastic representation $X = \mu + D\Sigma^{1/2}U$, where $D \perp U$ with $D^2 \sim \chi_d^2$ and U being uniformly distributed on the unit sphere.

Exercise 3.7.4 (Correlation test). Given a sample of size $n > 2$ from a bivariate normal distribution, use the same type of procedure as in Example 3.2.5 to derive an exact test of the hypothesis that the two variates are uncorrelated. That is: specify u and v , find a function of them that is parameter-free under H_0 , conclude independence, and go over to the marginal distribution. Finally transform to a test statistic of known distribution. Hint: $\sqrt{n-2} \frac{r}{\sqrt{1-r^2}}$ is exactly t_{n-2} -distributed under H_0 .

Exercise 3.7.5. Suppose we perform 10 tests (e.g. test the association between 10 different outcomes and a potential prognostic factor) and obtain the p -values 0.0140, 0.2960, 0.9530, 0.0031, 0.1050, 0.6410, 0.7810, 0.9010, 0.0053, 0.4500.

1. Which hypotheses are rejected after Bonferroni correction?
2. Which hypotheses are rejected after Holm correction?
3. Which hypotheses are rejected using the Benjamini and Hochberg procedure?

Exercise 3.7.6. Consider $2m$ independent coins $X_i \sim \text{Bern}(\theta_i)$ $i = 1, \dots, m$. Suppose that $\theta_i = \frac{1}{2} + \frac{i}{2m+1}$ for $i = 1, \dots, m$ and $\theta_i = \frac{1}{2}$ for $i = m+1, \dots, 2m$. For each of the coins we make n independent tosses.

- (a) For each i find the most powerful test for testing $\theta = \frac{1}{2}$ against $\theta > \frac{1}{2}$ for level $\alpha \leq 0.05$ (finding the test of size exactly α may be too hard).
- (b) Consider now the multiple testing problem for all m coins. We proved that Bonferroni and the Holm procedures both control the FWER. Describe how Bonferroni will look in this case. We want $\text{FWER} \leq 0.05$.
- (c) Take a look at the power of Bonferroni. Denote by \mathbb{P}_i the distribution $\text{Bern}(\theta_i)$ for $i = 1, \dots, m$ (these are the coins for which the null does not hold). Try to find some sufficient conditions to bound the probability of the type II error for each of the false nulls by, say, 0.05.
- (d) Provide some simulations to see the difference in power between the Bonferroni and Holm procedures. For fixed m consider $n = m, 10m, 100m$ to see how the answer depends on the ratio between n and m .

(some parts of this problem can be approached in various ways so there is no one correct solution)

Part II

**Statistical Learning Theory:
An Empirical Process
Perspective**

4

Motivation and examples (1 week)

This part of the lecture focuses on statistical learning theory taking the empirical process perspective. These techniques have become now a standard toolbox for studying modern statistical scenarios. We will introduce basic tools in empirical processes and apply those tools in mathematical statistics and machine learning. We will focus on the non-asymptotic perspective¹.

We start by introducing the main objects of this theory and some motivating examples².

4.1 Uniform law of large numbers

Suppose X, X_1, \dots, X_n are independent and identically distributed random variables taking values in $\mathcal{X} \subseteq \mathbb{R}$. Let the CDF of the underlying distribution be $F(t) = \mathbb{P}(X \leq t)$ for $t \in \mathbb{R}$. The empirical CDF \hat{F}_n is

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

As $\mathbb{E}\mathbb{1}\{X \leq t\} = \mathbb{P}(X \leq t)$, it is clear that, for every $t \in \mathbb{R}$, $\mathbb{E}\hat{F}_n(t) = F(t)$. By the strong law of large numbers, for every t , $\hat{F}_n(t) \xrightarrow{\text{a.s.}} F(t)$ as $n \rightarrow \infty$. The following stronger result is well-known.

Theorem 4.1.1 (Glivenko-Cantelli). *For any distribution, the empirical CDF \hat{F}_n is a strongly consistent estimator of the population CDF in the uniform norm, meaning that*

$$\|\hat{F}_n - F\|_\infty := \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \xrightarrow{\text{a.s.}} 0.$$

We will provide the proof of this result using a more general theory in Section 6.3; c.f. Proposition 6.3.6.

One reason, this result can be useful is because in many situations we want to estimate some functional $\gamma(F)$ of the population CDF. For example $\gamma_g(F) := \int g(x)dF(x)$ is the expectation $\mathbb{E}g(X)$. Also, for

¹ Good references:

Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019; and A. W. van der Vaart and Jon A. Wellner. *Weak convergence and empirical processes—with applications to statistics*. Springer Series in Statistics. Springer, Cham, 2023. Second edition

² Thanks go to Qiang Sun for inspiring this. Part of the material was discussed by Wenlong in the Fall semester.

any $\alpha \in [0, 1]$, the quantile functional Q_α is given

$$Q_\alpha(F) := \inf\{t \in \mathbb{R}; F(t) \geq \alpha\}.$$

If we have access to F only through the sample CFD \hat{F}_n . A natural way to estimate $\gamma(F)$ is using the plug-in estimator $\gamma(\hat{F}_n)$. We could define convergence $\hat{F}_n \rightarrow F$ in the sup-norm, as $\|\hat{F}_n - F\|_\infty \xrightarrow{\text{a.s.}} 0$. Then, we get that $\gamma(\hat{F}_n)$ is almost surely consistent for $\gamma(F)$ as long as the functional γ is continuous with respect to the sup-norm. Indeed, by continuity in the sup-norm

$$\{\omega : \|\hat{F}_n - F\|_\infty \rightarrow 0\} \subseteq \{\omega : \|\gamma(\hat{F}_n) - \gamma(F)\|_\infty \rightarrow 0\}.$$

By Theorem 4.1.1, the set on the left has measure 1 and so the set on the right also has measure one proving $\gamma(\hat{F}_n) \xrightarrow{\text{a.s.}} \gamma(F)$.

We are also interested in the following generalization of Theorem 4.1.1. Let \mathcal{F} denote a class of integrable real-valued functions on \mathcal{X} , and let $\{X_i\}_{i=1}^n$ be a collection of i.i.d. samples from some distribution \mathbb{P} over \mathcal{X} . We want to analyze the following quantity

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right|, \quad (4.1)$$

which measures the absolute deviation between the sample average $\frac{1}{n} \sum f(X_i)$ and the population average $\mathbb{E}[f(X)]$, uniformly over class \mathcal{F} .

Definition 4.1.2. We say that \mathcal{F} is a **Glivenko-Cantelli class** for \mathbb{P} if $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \rightarrow 0$ almost surely as $n \rightarrow \infty$. If the convergence in probability holds, we say that \mathcal{F} is a **weak Glivenko-Cantelli class** for \mathbb{P} .

If \mathcal{F} is the family of indicator functions $f(x) = \mathbb{1}\{x \leq t\}$ for some $t \in \mathbb{R}$, we recover the CDF example in Theorem 4.1.1. In Chapter 6 we discuss some other examples of **Glivenko-Cantelli classes**. In general, the condition is that the class cannot be too rich as illustrated by the example below.

Example 4.1.3. Suppose that \mathcal{F} is the set of all indicator functions $\mathbb{1}_A(x)$ for all measurable sets $A \subseteq \mathbb{R}$ and suppose that \mathbb{P} is absolutely continuous with respect to the Lebesgue measure on \mathbb{R} . We have $\mathbb{E}\mathbb{1}_A(X) = \mathbb{P}(X \in A)$. For any $x_1, \dots, x_n \in \mathbb{R}$ the set $A = \mathbb{R} \setminus \{x_1, \dots, x_n\}$ is measurable and $\mathbb{P}(X \in A) = 1$. We thus get that $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} = 1$ for all n , and so, this class of indicator functions is not Glivenko-Cantelli.

Even if \mathcal{F} is not a Glivenko-Cantelli class, we still want to control the size (4.1) by answering the following questions:

1. Does the random variable $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ in (4.1) concentrate around its expectation? What is the size of this expectation?

Recall that a function $f : \mathcal{X} \rightarrow \mathbb{R}$ on a metric space (\mathcal{X}, d) is continuous if $p_n \rightarrow p$ in \mathcal{X} implies that $f(p_n) \rightarrow f(p)$ in \mathbb{R} . Also recall that $p_n \rightarrow p$ is equivalent to $d(p_n, p) \rightarrow 0$.

Note that there may be measurability concerns associated with this definition. We will skip the details; see Section 4.4 in Wainwright's book for some details.

2. Provide finite-sample bounds on $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$, that is, bounds that hold for every n , in terms of the class function and the distribution \mathbb{P} of X .

The strategy for controlling (4.1) is to first show that this quantity concentrates around its mean. Given the techniques described in Chapter 5, this is often easy³. Second, we need to control the mean. This is often more complicated and requires more advanced techniques that we overview in Chapter 6. Before we discuss more technical aspects, we first discuss one motivating example in machine learning and one in statistics.

4.1.1 Example: Binary classification

Consider a random vector (X, Y) with values in $\mathcal{X} \times \{-1, 1\}$. A classifier is a function $g : \mathcal{X} \rightarrow \{-1, 1\}$. The error of the classifier is given by

$$R(g) := \mathbb{P}(g(X) \neq Y).$$

The goal of binary classification is to construct a classifier with small error based on n *i.i.d.* observations $(X_1, Y_1), \dots, (X_n, Y_n)$ having the same distribution as (X, Y) . The empirical error of the classifier g is

$$R_n(g) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(X_i) \neq Y_i).$$

A natural strategy for classification is to pick a class of classifiers \mathcal{C} and then to choose the classifier in \mathcal{C} which has the smallest **training error**

$$\hat{g}_n := \arg \min_{g \in \mathcal{C}} R_n(g).$$

Here one possible choice is to use the logistic regression in Example 1.7.2. A good classifier should have a small **out of sample test error**

$$R(\hat{g}_n) = \mathbb{P}(\hat{g}_n(X) \neq Y | (X_1, Y_1), \dots, (X_n, Y_n)).$$

One possible question concerns how close is \hat{g}_n to the optimal classifier $g^* = \arg \min_{g \in \mathcal{C}} R(g)$. To compare $R(\hat{g}_n)$ with $R(g^*)$ note that

$$\begin{aligned} R(\hat{g}_n) &= R(g^*) + R(\hat{g}_n) - R_n(\hat{g}_n) + R_n(\hat{g}_n) - R(g^*) \\ &\leq R(g^*) + R(\hat{g}_n) - R_n(\hat{g}_n) + R_n(g^*) - R(g^*) \quad (4.2) \\ &\leq R(g^*) + 2 \sup_{g \in \mathcal{C}} |R_n(g) - R(g)|. \end{aligned}$$

Another question could be about comparing the training error and the test error, which amounts to comparing $R(\hat{g}_n)$ with $R_n(\hat{g}_n)$. Here

³ For instance, for the class in Example 4.1.3, we show concentration around the mean in Example 5.3.6.

In the language of the statistical decision theory we first define the loss $L(g) = \mathbb{1}\{g(X) \neq Y\}$. Then the error $R(g)$ is simply the risk of the classifier g .

← Exercise 4.3.1

It is natural to associate each classifier with the set $\{x : g(x) = 1\}$. Thus, the class \mathcal{C} can be identified with a set of measurable subsets of \mathcal{X} .

we have

$$R(\widehat{g}_n) = R_n(\widehat{g}_n) + R(\widehat{g}_n) - R_n(\widehat{g}_n) \quad (4.3)$$

$$\leq R_n(\widehat{g}_n) + \sup_{g \in \mathcal{C}} |R_n(g) - R(g)|. \quad (4.4)$$

With the analogous bound on $R_n(\widehat{g}_n)$ we conclude that $|R_n(\widehat{g}_n) - R(\widehat{g}_n)| \leq \sup_{g \in \mathcal{C}} |R_n(g) - R(g)|$. The key quantity in both cases is

$$\sup_{g \in \mathcal{C}} |R_n(g) - R(g)|,$$

which is a special case of (4.1), when \mathcal{F} is taken to be the class of all functions $\mathbb{1}(g(x) \neq y)$ as g varies over \mathcal{C} , where the data are (X_i, Y_i) instead of X_i . We provide more details in Section 6.3.1 after developing necessary theory.

4.1.2 Example: M-Estimation

Let X, X_1, \dots, X_n are i.i.d. from \mathbb{P} where $\mathbb{P} \in \mathcal{P}$, with $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$ for $\Theta \subseteq \mathbb{R}^d$. A popular method of finding an estimator $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ is to minimize a criterion function of the form

$$\theta \mapsto M_n(\theta) = \frac{1}{n} \sum_{i=1}^n m_\theta(X_i).$$

Here $m_\theta : \mathcal{X} \mapsto \mathbb{R} \cup \{+\infty\}$ are known functions. An estimator maximizing $M_n(\theta)$ is called an M-estimator. Often the minimizer is sought by setting a derivative equal to zero. Therefore, the name M-estimator is also used for estimators satisfying systems of equations of the type

$$\Psi_n(\theta) = \frac{1}{n} \sum_{i=1}^n \psi_\theta(X_i) = \mathbf{0}. \quad (4.5)$$

Here $\psi_\theta : \mathcal{X} \rightarrow \mathbb{R}^d$ are known vector-valued maps.

Example 4.1.4 (Maximum likelihood estimators). *Suppose X, X_1, \dots, X_n have a common density p_θ . Then the maximum likelihood estimator maximizes the likelihood $\prod_i p_\theta(X_i)$, or equivalently*

$$\ell_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i).$$

This is an expectation of $\log p(X)$ with respect to the sample distribution.

Thus a maximum likelihood estimator is an M-estimator with $m_\theta = -\log p_\theta$. If the density is differentiable with respect to θ for each fixed x , then the maximum likelihood estimator also solves an equation of type (4.5) with $\psi_\theta = \nabla_\theta \log p_\theta$, the score function of the model.

Other simple examples of M-estimators include the sample median (c.f. Exercise 2.7.2), the least squares estimator, and more generally, the estimators obtained by minimizing the empirical risk (2.1) in

which case $m_\theta(x) = L(\theta, \delta(x))$. The general theory of M-estimators originated in robust statistics.

Example 4.1.5 (Huber Robust loss). Consider a regression problem. Let $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the design matrix with rows $\mathbf{x}_1, \dots, \mathbf{x}_n$. Let $L(y, a) = \frac{1}{2}(y - a)^2$. In the standard least squares approach, the vector of coefficient is estimated by minimizing the empirical risk

$$m(\beta) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|^2 = \frac{1}{n} \sum_{i=1}^n L(y_i, \beta^T \mathbf{x}_i).$$

Consider instead the loss

$$L_\delta(y, a) = \begin{cases} \frac{1}{2}(y - a)^2 & \text{if } |y - a| \leq \delta \\ \delta(|y - a| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}.$$

Note that $L_\delta(y, a)$ is differentiable everywhere.

The corresponding M-estimator is the minimizer of

$$m_\delta(\beta) = \frac{1}{n} \sum_{i=1}^n L_\delta(y_i, \beta^T \mathbf{x}_i).$$

Let $M(\theta) = \mathbb{E}m_\theta(X)$. In the theory of M-estimation, the target quantity for the estimator $\hat{\theta}_n$ is

$$\theta^* := \arg \min_{\theta \in \Theta} M(\theta).$$

For example, suppose data come from the distribution \mathbb{P} with density q . Let $m_\theta(X) = -\log p_\theta(X)$. Then

$$\arg \min_{\theta} -\mathbb{E} \log p_\theta(X) = \arg \min_{\theta} \mathbb{E} \log \frac{q(X)}{p_\theta(X)},$$

which is simply the Kullback-Leibler divergence between the distribution \mathbb{P} and \mathbb{P}_θ . In particular, if $\mathbb{P} \in \mathcal{P}$ then it is the unique minimizer.

The main question of interest while studying M-estimators concerns accuracy of $\hat{\theta}_n$ for estimating θ^* . In the asymptotic regime $n \rightarrow \infty$, the two key questions are:

1. Is $\hat{\theta}_n$ consistent for estimating θ^* , equivalently, does $d(\hat{\theta}_n, \theta^*) \xrightarrow{P} 0$; see Section 8.1 for relevant definitions and basic results.
2. If yes, what is the rate of convergence of $d(\hat{\theta}_n, \theta^*)$ to zero? The usual rate of convergence is $O_p(n^{-1/2})$, so that $\sqrt{n}(\hat{\theta}_n - \theta^*)$ converges in distribution⁴.
3. How shall we do inference on θ^* using $\hat{\theta}_n$. For this we may need to understand better how $\hat{\theta}_n$ concentrates around θ^* .

⁴ For a simple example see Exercise 8.6.5.

Some basic asymptotic results for M-estimators are provided in Appendix 8.3. Here, we complement this with the link to the uniform law of large numbers. For the first question we investigate closeness of $M_n(\theta) = \frac{1}{n} \sum_{i=1}^n m_\theta(X_i)$ to $M(\theta) = \mathbb{E}m_\theta(X)$ in some sort of uniform sense over θ , which leads to investigation of (4.1) for

$$\mathcal{F} = \{m_\theta : \theta \in \Theta\}$$

and then translate the result back to the result on $d(\hat{\theta}_n, \theta^*)$.

This is how this can be done: Let $D_\epsilon = \{\theta \in \Theta : d(\theta, \theta^*) \geq \epsilon\}$. We can bound $d(\hat{\theta}_n, \theta^*)$ as

$$\begin{aligned} \mathbb{P}(d(\hat{\theta}_n, \theta^*) \geq \epsilon) &\leq \mathbb{P}(\sup_{\theta \in D_\epsilon} (M_n(\theta^*) - M_n(\theta)) \geq 0) \\ &\leq \mathbb{P}\left(\sup_{\theta \in D_\epsilon} \{(M_n(\theta^*) - M(\theta^*)) - (M_n(\theta) - M(\theta))\} \geq -\sup_{\theta \in D_\epsilon} \{M(\theta^*) - M(\theta)\}\right) \\ &\leq \mathbb{P}\left(2 \sup_{\theta \in D_\epsilon} \{|M_n(\theta) - M(\theta)|\} \geq \inf_{\theta \in D_\epsilon} \{M(\theta) - M(\theta^*)\}\right), \end{aligned}$$

where going from the first to the second line we used the fact that $\sup_\theta f(\theta) \leq \sup_\theta g(\theta) + \sup_\theta (f(\theta) - g(\theta))$. Note that the bound above again relates to the random quantity (4.1). If θ^* is uniquely and globally identifiable, that is, there exists $\eta > 0$ depending on ϵ such that $\inf_{\theta \in D_\epsilon} (M(\theta) - M(\theta^*)) > \eta$. Thus, with proper control of the size of $\sup_{\theta \in D_\epsilon} |M_n(\theta) - M(\theta)|$, the right-hand side in the above display diminishes to zero⁵.

⁵ We get $\mathbb{P}(d(\hat{\theta}_n, \theta^*) \geq \epsilon) \leq \mathbb{P}(\sup_{\theta \in D_\epsilon} |M_n(\theta) - M(\theta)| \geq \eta/2)$.

4.2 The Uniform Central Limit Theory

The classical central limit theorem (CLT) studies the following type of results: for an *i.i.d.* sequence X, X_1, \dots, X_n

$$\mathbf{G}_n(f) := \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right) \rightsquigarrow N(0, \text{var}(f(X))). \quad (4.6)$$

The uniform central limit theorem studies the above convergence in distribution uniformly over f in the class \mathcal{F} . To illustrate this idea consider the following example.

4.2.1 Example: Uniform empirical process

Suppose X, X_1, \dots, X_n are *i.i.d.* uniform on $[0, 1]$. Let $\mathcal{F} = \{\mathbb{1}_{(-\infty, t]}(x) : t \in \mathbb{R}\}$. Define

$$U_n(t) = \sqrt{n}(\hat{F}_n(t) - F(t)), \quad t \in \mathbb{R}, \quad (4.7)$$

where $\hat{F}_n(t) = \frac{1}{n} \sum_i \mathbb{1}\{X_i \geq t\}$ and $F(t) = \mathbb{P}(X \leq t) = t$.

This defines a stochastic process $\{U_n(t) : t \in \mathbb{R}\}$ which is typically referred to as an **empirical process** or, more specifically in this case, a uniform empirical process.

How do we find a candidate for the limit? The CLT states that, for each $t \in [0, 1]$, $U_n(t) \rightsquigarrow N(0, t(1-t))$ as $n \rightarrow \infty$ ⁶. Moreover, for every fixed t_1, \dots, t_k , the multivariate CLT states that

$$(U_n(t_1), \dots, U_n(t_k)) \rightsquigarrow N(0, \Sigma),$$

where $\Sigma = (\Sigma_{ij})$, $\Sigma_{ij} = t_i \wedge t_j - t_i t_j$. This follows simply because we can rewrite $U_n(t) = \frac{1}{\sqrt{n}} \sum_i (\mathbb{1}(X_i \leq t) - t)$ and, for any t, t' ,

$$\begin{aligned} \mathbb{E}U_n(t)U_n(t') &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\mathbb{1}(X_i \leq t) - t)(\mathbb{1}(X_i \leq t') - t')] \\ &= \frac{1}{n} \sum_{i=1}^n (t \wedge t' - tt') = t \wedge t' - tt'. \end{aligned}$$

In consequence, the candidate for the limiting process is the Gaussian process with kernel $\kappa(s, t) = s \wedge t - st$.

We will be interested in the underlying limiting process called a Brownian bridge.

Definition 4.2.1 (Brownian bridge). *Brownian bridge* $\{U(t) : 0 \leq t \leq 1\}$ is a stochastic process satisfying the following conditions

1. Every realization is continuous in $[0, 1]$ with $U(0) = U(1) = 0$.
2. For every fixed t_1, \dots, t_k , the vector $(U(t_1), \dots, U(t_k)) \sim N(0, \Sigma)$, where $\sigma_{ij} = t_i \wedge t_j - t_i t_j$.

Thus, the finite dimensional representation of $\{U_n(t) : t \in [0, 1]\}$ converges in distribution to that of $\{U(t) : t \in [0, 1]\}$. It is then natural to ask whether the entire process $\{U_n(t)\}$ converges in distribution to the process $\{U(t)\}$.

Convergence of stochastic processes can be defined as follows. First recall the definition of convergence in distribution given in Section 8.1. Equivalently, by Portmanteau Lemma 8.1.1, we say that a random sequence (Z_n) with values in \mathbb{R}^k converges in distribution to Z if and only if $\mathbb{E}h(Z_n) \rightarrow \mathbb{E}h(Z)$ for every $h \in \mathcal{C}_b(\mathbb{R}^k)$, where $\mathcal{C}_b(\mathbb{R}^k)$ denotes the class of real-valued bounded functions on \mathbb{R}^k . This equivalent definition can now be generalized to stochastic processes.

Definition 4.2.2. For any set \mathcal{F} by $\ell^\infty(\mathcal{F})$ denote the set of bounded real-valued functions on \mathcal{F} with the sup-norm.

In particular, $\ell^\infty([0, 1])$ is the space of all bounded functions on $[0, 1]$ with metric $\rho(f, g) := \sup_{t \in [0, 1]} |f(t) - g(t)|$.

⁶ We have $\mathbb{E}U_n(t) = 0$ and $\text{var}(U_n(t)) = \text{var}(\mathbb{1}\{X \leq t\})$

A Brownian bridge is an example of a Gaussian process on $[0, 1]$ with the underlying kernel function $\kappa(s, t) = s \wedge t - st$.

Definition 4.2.3. The process $\{U_n(t)\}$ converges to $\{U(t)\}$ in distribution if

$$\mathbb{E}h(U_n) \rightarrow \mathbb{E}h(U) \quad \text{as } n \rightarrow \infty$$

for every bounded and continuous real-valued function $h : \ell^\infty[0, 1] \rightarrow \mathbb{R}$.

There is one measure-theoretic issue with this definition as for some h , $h(U_n)$ could not be measurable. A technical remedy is to replace expectation by the outer expectation $\mathbb{E}^*h(U_n)$ induced by the outer measure. For details see Chapter 1 in ⁷.

4.2.2 The general case

Consider the general empirical process $\mathbb{G}_n(f)$ defined in (4.6), where $f \in \mathcal{F}$. Under the assumption that $\sup_{f \in \mathcal{F}} |f(x)| < \infty$ for every $x \in \mathcal{X}$, the function $f \mapsto \mathbb{G}_n(f)$ belongs to $\ell^\infty(\mathcal{F})$.

Definition 4.2.4. We say that \mathcal{F} is a **Donsker class** (or \mathbb{P} -Donsker) if $\mathbb{G}_n(\mathcal{F}) = \{\mathbb{G}_n(f) : f \in \mathcal{F}\}$ converges in distribution in $\ell^\infty(\mathcal{F})$ to some $\mathbb{G}(\mathcal{F}) = \{\mathbb{G}(f) : f \in \mathcal{F}\}$ as $n \rightarrow \infty$.

The limiting process $\mathbb{G}(\mathcal{F})$ is a Gaussian process: for every f_1, \dots, f_k the vector $(\mathbb{G}(f_1), \dots, \mathbb{G}(f_k))$ is a multivariate Gaussian distribution.

This all looks abstract but the following example should illustrate importance of these considerations in statistics.

Example 4.2.5 (A goodness-of-fit statistics). Consider X, X_1, \dots, X_n i.i.d. from a distribution \mathbb{P} on \mathbb{R} with CDF F . Suppose we want to test $H_0 : F = F_0$ versus $H_1 : F \neq F_0$. Kolmogorov proposed the following test statistics

$$D_n := \sqrt{n} \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F_0(x)|. \quad (4.8)$$

It turns out that, under H_0 , the distribution of D_n does not depend on F_0 ; see Exercise 4.3.2. If F_0 is strictly increasing and continuous we can then with no loss of generality assume that F_0 is a uniform distribution on $[0, 1]$; $F_0(t) = t$ for $t \in [0, 1]$. We obtain that under the null

$$D_n = \sup_{0 \leq t \leq 1} |U_n(t)| = \|U_n\|_\infty \rightarrow \|U\|_\infty,$$

where $U_n(t)$ was defined in (4.7). Thus⁸

$$\lim_{n \rightarrow \infty} \mathbb{P}(D_n \leq x) = \mathbb{P}\left(\sup_{0 \leq t \leq 1} |U(t)| \leq x\right).$$

Example 4.2.6 (Asymptotics of MLE). Suppose X, X_1, \dots, X_n are iid from \mathbb{P}_{θ_0} . The maximum likelihood estimator is

$$\hat{\theta}_n := \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta).$$

⁷ A. W. van der Vaart and Jon A. Wellner. *Weak convergence and empirical processes—with applications to statistics*. Springer Series in Statistics. Springer, Cham, 2023. Second edition

← Exercise 4.3.2

⁸ $\mathbb{P}_n(D_n \leq x) = \mathbb{E}\mathbb{1}\{D_n \leq x\} = \mathbb{E}\mathbb{1}\{\|U_n\|_\infty \leq x\}$

A classical result is that, under second order or third order smoothness conditions $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, \mathbf{I}(\theta_0)^{-1})$, where $\mathbf{I}(\theta_0)$ is the Fisher information matrix defined as

$$\mathbf{I}(\theta_0) = \mathbb{E}[\nabla \log p_\theta(X) \nabla^\top \log p_\theta(X)] \Big|_{\theta=\theta_0}.$$

See Section 8.3.2 for some details.

What is the minimal smoothness assumption that is needed for the above result to hold? It turns out that if we use UCLT together with the notion of differentiability in quadratic mean (DQM), first-order smoothness will be enough. This is useful in some examples that involve the Laplace density.

Example 4.2.7 (Asymptotics of M-estimators). The UCLT can also be used to derive asymptotic distributions for M-estimators. We consider two representative examples of location estimators.

1. The sample median is defined as

$$\hat{\theta}_n = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n |X_i - \theta|.$$

Assume the CDF F is differentiable around its median θ_0 with positive derivative $F'(\theta_0) =: p(\theta_0)$. Using UCLT, we can prove that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, \frac{1}{4p^2(\theta_0)}).$$

2. A form of mode estimator can be defined as

$$\hat{\theta}_n = \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n m_\theta(X_i),$$

where $m_\theta(x) = \mathbb{1}(|x - \theta| \leq 1)$. For this estimator, the asymptotic distribution is very complicated as we have

$$n^{1/3}(\hat{\theta}_n - \theta) \rightsquigarrow \arg \max_{h \in \mathbb{R}} \{aZ(h) - bh^2\},$$

where Z is a standard two-sided Brownian motion starting from 0, and

$$a^2 = p(\theta_0 + 1) - p(\theta_0 - 1) \quad \text{and} \quad b = \frac{1}{2}(p'(\theta_0 + 1) - p'(\theta_0 - 1)),$$

where p is the density function, unimodal, and symmetric.

This is θ that contains the most of X_i 's in its unit neighborhood.

4.3 Exercises

Exercise 4.3.1. In the binary classification problem in Section 4.1.1, show that the Bayes classifier g_0 given by

$$g_0(x) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1|X = x) \geq \frac{1}{2}, \\ -1 & \text{otherwise} \end{cases}$$

minimizes the misclassification probability $R(g)$ over all measurable functions g .

Exercise 4.3.2. Show that the the statistics (4.8) does not depend on F_0 .
Hint: First assume that the CDF is continuous and strictly increasing. For the general result use (3.24).

5

Concentration of measure (3 weeks)

In this chapter we are interested in bounding random fluctuations of functions of many independent random variables. Variables X, X_1, \dots, X_n are independent and take values in some \mathcal{X} . Let $g : \mathcal{X}^n \rightarrow \mathbb{R}$ and

$$Z = g(X_1, \dots, X_n).$$

The function g can be quite complex. For example, as motivated in the previous chapter, we could have

$$Z = \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right|.$$

How large are “typical” deviations of Z from $\mathbb{E}Z$? In particular, we seek upper bounds for

$$\mathbb{P}(Z \geq \mathbb{E}Z + t) \quad \text{and} \quad \mathbb{P}(Z \leq \mathbb{E}Z - t)$$

for $t > 0$. When $\mathbb{E}Z$ is unknown, to obtain direct bounds on the concentration of Z , explicit bounds on $\mathbb{E}Z$ will also be needed, which will lead to more advanced considerations in Chapter 6.

There are various methods that include: martingales, information theoretic and transportation methods, Talagrand’s induction method, and logarithmic Sobolev inequalities. We will not get into too many details here. We refer to two excellent books ¹ and ² for a thorough overview of the theory.

5.1 Basic inequalities

In this section we set up the scene for more advanced considerations by recalling some basic probability inequalities. The simplest inequality is the Markov’s inequality.

Proposition 5.1.1 (Markov’s inequality). *If $Z \geq 0$ and $t > 0$ then*

$$\mathbb{P}(Z \geq t) \leq \frac{\mathbb{E}Z}{t}.$$

¹ Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013

² Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019

Proof. We have

$$t\mathbb{P}(Z \geq t) = \mathbb{E}[t\mathbb{1}(Z \geq t)] \leq \mathbb{E}[Z\mathbb{1}(Z \geq t)] \leq \mathbb{E}Z,$$

where the last inequality follows because $Z \geq 0$. \square

Markov's inequality with its elementary proof looks very innocent. It is then surprising to see how many powerful results can be obtained from it. We first show that, for distributions for which the second moment exists, the Markov's inequality implies the Chebyshev's inequality.

Proposition 5.1.2 (Chebyshev's inequality). *For every $t > 0$ we have*

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) = \mathbb{P}((Z - \mathbb{E}Z)^2 \geq t^2) \leq \frac{\text{var}(Z)}{t^2}.$$

A particularly important instance is when Z is a sum of i.i.d. random variables, $Z = \sum_{i=1}^n X_i$. Denote $\mu = \mathbb{E}X$, $\sigma^2 = \text{var}(X)$. In this case $\text{var}(Z) = n\sigma^2$ and we get

$$\mathbb{P}\left(\left|\sum_{i=1}^n X_i - n\mu\right| \geq t\right) \leq \frac{n\sigma^2}{t^2}$$

or equivalently, denoting $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$,

$$\mathbb{P}(|\sqrt{n}(\bar{X}_n - \mu)| \geq t) \leq \frac{\sigma^2}{t^2}.$$

This bound is however not very tight. Let $U \sim N(0, 1)$ and note that, by the Central Limit Theorem,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\sqrt{n}(\bar{X}_n - \mu)| \geq t) = 2\mathbb{P}(U \geq \frac{t}{\sigma}) \leq \frac{\sigma}{t} e^{-t^2/(2\sigma^2)},$$

where the last inequality is part of Exercise 5.5.2. In particular, at least for very large n , we expect an exponential decrease in t^2/σ^2 .

The trick to get the expected rate of decrease is to use the Markov's inequality in a more clever way.

Proposition 5.1.3 (Chernoff bounds). *Suppose that the moment generating function of Z exists in some neighborhood $(-b, b)$ of zero. Then for every $\lambda > 0$ in this neighborhood*

$$\mathbb{P}(Z - \mathbb{E}Z \geq t) = \mathbb{P}(e^{\lambda(Z - \mathbb{E}Z)} \geq e^{\lambda t}) \leq \frac{\mathbb{E}e^{\lambda(Z - \mathbb{E}Z)}}{e^{\lambda t}}.$$

In consequence,

$$\log \mathbb{P}(Z - \mathbb{E}Z \geq t) \leq \inf_{\lambda \in [0, b)} \{\log \mathbb{E}e^{\lambda(Z - \mathbb{E}Z)} - \lambda t\}.$$

More generally, if the corresponding k -th moment exists for some $k > 0$ then

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq \frac{\mathbb{E}(|Z - \mathbb{E}Z|^k)}{t^k}$$

← Exercise 5.5.2

Note that if $M(\lambda) = \mathbb{E}e^{\lambda(Z-\mathbb{E}Z)}$ is the moment generating function of $Z - \mathbb{E}Z$ then $K(\lambda) = \log M(\lambda)$ is the cumulant generating function. We get

$$\log \mathbb{E}e^{\lambda(Z-\mathbb{E}Z)} - \lambda t = K(\lambda) - \lambda t.$$

By Hölder's inequality, $K(\lambda)$ is a convex function with the same proof as in Theorem 1.3.2. We will later discuss some techniques to bound the moment generating function.

A useful observation is that, if X, X_1, X_2, \dots, X_n are i.i.d. then

$$\mathbb{P}(\bar{X}_n - \mu \geq t) \leq \frac{\mathbb{E}e^{\lambda(\bar{X}_n - \mu)}}{e^{\lambda t}} = \left(\frac{\mathbb{E}e^{\frac{\lambda}{n}(X - \mathbb{E}X)}}{e^{\frac{\lambda}{n}t}} \right)^n.$$

Optimizing over λ is equivalent to optimizing over λ/n and so a Chernoff bound on a single X , directly gives a Chernoff bound for an average of its independent copies.

Example 5.1.4. Let Y, Y_1, \dots, Y_n be an i.i.d. sample such that $Y = X^2$ with $X \sim N(0, 1)$. First note that

$$\mathbb{E}e^{\lambda Y} = \begin{cases} \sqrt{\frac{1}{1-2\lambda}} & \text{for } \lambda < \frac{1}{2}, \\ +\infty & \text{otherwise.} \end{cases}$$

Thus, for $\lambda < 1/2$,

$$\mathbb{P}(Y - 1 \geq t) \leq \sqrt{\frac{1}{1-2\lambda}} \frac{1}{e^{\lambda(t+1)}}.$$

The optimal $\lambda^* = \frac{1}{2} \frac{t}{t+1}$, which gives the Chernoff bound³

$$\mathbb{P}(Y - 1 \geq t) \leq \sqrt{t+1} e^{-t/2}$$

and consequently

$$\mathbb{P}(\bar{Y}_n - 1 \geq t) \leq (t+1)^{n/2} e^{-nt/2}.$$

It is convenient to have an explicit exponential bound for small t . In this case we can use the fact that $\log(1+t) - t \leq -t^2/4$ for $t \in [0, 1]$ ⁴ to conclude that, for every $t \in [0, 1]$,

$$\mathbb{P}(\bar{Y}_n - 1 \geq t) \leq (t+1)^{n/2} e^{-nt/2} = e^{\frac{n}{2}(\log(t+1)-t)} \leq e^{-nt^2/8}.$$

A simple yet powerful illustration of the Chernoff bounds in Proposition 5.1.3 is given by following example.

Example 5.1.5 (Johnson-Lindenstrauss lemma). Let $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^D$ and $\epsilon > 0$. We are looking for a function $f : \mathbb{R}^D \rightarrow \mathbb{R}^d$ with $d < D$ such that

$$(1 - \epsilon) \|\mathbf{a}_i - \mathbf{a}_j\|^2 < \|f(\mathbf{a}_i) - f(\mathbf{a}_j)\|^2 < (1 + \epsilon) \|\mathbf{a}_i - \mathbf{a}_j\|^2 \quad (5.1)$$

The convex conjugate to $K(\lambda)$ is $K^*(t) = \inf_{\lambda \in \mathbb{R}} \{K(\lambda) - \lambda t\}$ and we get $\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq e^{K^*(t)}$.
← Exercise 5.5.3

³ Plot this to see that this bound is actually not that great for $t \in (0, 5)$.

⁴ To show that $-t^2/4 - \log(1+t) + t \geq 0$ for $t \in [0, 1]$ we note that the left hand side is zero for $t = 0$ and its derivative is nonnegative for $t \in [0, 1]$.

← Exercise 5.5.4

for all $i, j = 1, \dots, n$. Johnson-Lindenstrauss lemma states such an embedding f exists if

$$d > \frac{16}{\epsilon^2} \log n. \quad (5.2)$$

(Note that the bound does not depend on D !) We present a probabilistic proof of this fact and show that f can be taken as a linear function. In fact, if we choose f at random, it works with positive probability. Let $W \in \mathbb{R}^{d \times D}$ be a random matrix such that the W_{ij} are independent $N(0, 1/d)$.

Take $f(\mathbf{a}) = W\mathbf{a}$. Then, denoting $\mathbf{b}_{ij} := (\mathbf{a}_i - \mathbf{a}_j) / \|\mathbf{a}_i - \mathbf{a}_j\|$, we can rewrite (5.1) as

$$1 - \epsilon < \|W\mathbf{b}_{ij}\|^2 < 1 + \epsilon.$$

or equivalently

$$\max_{i,j} \left| \|W\mathbf{b}_{ij}\|^2 - 1 \right| < \epsilon. \quad (5.3)$$

We first show that $\mathbb{E}\|W\mathbf{b}_{ij}\|^2 = 1$. Indeed, for any $\mathbf{b} \in \mathbb{R}^D$,

$$\mathbb{E}\|W\mathbf{b}\|^2 = \mathbb{E}(\mathbf{b}^\top W^\top W \mathbf{b}) = \mathbf{b}^\top \mathbb{E}(W^\top W) \mathbf{b} = \mathbf{b}^\top \mathbf{b} = \|\mathbf{b}\|^2. \quad (5.4)$$

Since $\|\mathbf{b}_{ij}\| = 1$ we get $\mathbb{E}\|W\mathbf{b}_{ij}\|^2 = 1$. To show (5.3) with positive probability, it is equivalent to show that

$$\mathbb{P}(\max_{i,j} \left| \|W\mathbf{b}_{ij}\|^2 - 1 \right| \geq \epsilon) < 1.$$

We next show that $\|W\mathbf{b}_{ij}\|^2$ can be written as $\frac{1}{d} \sum_{i=1}^d Z_i^2$ with $Z_i \sim N(0, 1)$. Indeed, for every $\mathbf{b} \in \mathbb{R}^D$ such that $\|\mathbf{b}\| = 1$

$$\|W\mathbf{b}\|^2 = (W\mathbf{b})^\top (W\mathbf{b}) = \sum_{i=1}^d (\mathbf{w}_i^\top \mathbf{b})^2,$$

where $\mathbf{w}_1, \dots, \mathbf{w}_d$ are the rows of W . Note that $\mathbb{E}(\mathbf{w}_i^\top \mathbf{b}) = 0$ and

$$\mathbb{E}((\mathbf{w}_i^\top \mathbf{b})^2) = \mathbf{b}^\top \mathbb{E}(\mathbf{w}_i \mathbf{w}_i^\top) \mathbf{b} = \frac{1}{d} \mathbf{b}^\top \mathbf{b} = \frac{1}{d}.$$

Thus, $\mathbf{w}_i^\top \mathbf{b}$ for $i = 1, \dots, d$ are i.i.d. $N(0, 1/d)$, we get

$$\|W\mathbf{b}\|^2 - 1 = \frac{1}{d} \sum_{i=1}^d (Z_i^2 - 1),$$

where Z_i are i.i.d. $N(0, 1)$. By the calculations in Example 5.1.4 and by (5.10), for every $\epsilon \leq 1$ we conclude that

$$\mathbb{P}(\left| \|W\mathbf{b}\|^2 - 1 \right| \geq \epsilon) \leq 2e^{-\epsilon^2 d/8}.$$

By the union bound

$$\mathbb{P}(\max_{i,j} \left| \|W\mathbf{b}_{ij}\|^2 - 1 \right| > \epsilon) \leq \binom{n}{2} \mathbb{P}(\left| \|W\mathbf{b}\|^2 - 1 \right| > \epsilon) \leq \binom{n}{2} 2e^{-\epsilon^2 d/8}.$$

It remains to check that if (5.2) holds then the right hand side is smaller than 1. But this is clear. Rewrite (5.2) as $de^2/8 > \log n^2$ and get

$$\binom{n}{2} 2e^{-e^2 d/8} = e^{\log(n(n-1)) - e^2 d/8} < e^{\log(n(n-1)) - \log n^2} = \frac{n-1}{n} < 1.$$

Note that the proof is not constructive, as the linear map satisfying (5.1) is not given explicitly.

5.2 Sub-gaussian and sub-exponential random variables

It is clear that if the values of Z concentrate, its variance must be small. To establish concentration using the Chebyshev's inequality in Proposition 5.1.2 it is important to obtain good bounds on the variance of Z . Similarly, in order to utilize the Chernoff bounds in Proposition 5.1.3 we need to obtain good bounds for the moment generating function. In this and in the next section we will focus on this last task.

To motivate the next definition note that if $X \sim N(\mu, \sigma^2)$ then

$$\mathbb{E}e^{\lambda(X-\mu)} = e^{\sigma^2 \lambda^2/2} \quad \text{for all } \lambda \in \mathbb{R}.$$

Substituting this into the Chernoff bound in Proposition 5.1.3 we get

$$\inf_{\lambda \geq 0} \left\{ \log \mathbb{E}e^{\lambda(X-\mu)} - \lambda t \right\} = \inf_{\lambda \geq 0} \left\{ \frac{\lambda^2 \sigma^2}{2} - \lambda t \right\} = -\frac{t^2}{2\sigma^2},$$

which gives

$$\mathbb{P}(X - \mu \geq t) \leq e^{-t^2/2\sigma^2} \quad \text{for all } t \geq 0. \quad (5.5)$$

In other words, for $\delta \in (0, 1)$,

$$\mathbb{P}\left(X - \mu \geq \sigma \sqrt{2 \log\left(\frac{1}{\delta}\right)}\right) \leq \delta.$$

The same concentration bounds hold in much greater generality for so called sub-Gaussian variables.

← Exercise 5.5.5

Definition 5.2.1. A variable X is sub-Gaussian (or σ -sub-Gaussian) if there is $\sigma > 0$ such that

$$\mathbb{E}e^{\lambda(X-\mu)} \leq e^{\sigma^2 \lambda^2/2} \quad \text{for all } \lambda \in \mathbb{R}$$

or equivalently

$$K_X(\lambda) - \lambda\mu \leq \sigma^2 \lambda^2/2 \quad \text{for all } \lambda \in \mathbb{R}, \quad (5.6)$$

where K_X is the cumulant generating function of X .

The following proposition shows that, if X takes values in a bounded interval $[a, b]$, then X is sub-Gaussian with $\sigma = (b - a)/2$.

Proposition 5.2.2. *Suppose that X has values in some bounded interval $[a, b]$. Then $\mathbb{E}e^{\lambda(X-\mathbb{E}X)} \leq e^{\lambda^2(b-a)^2/8}$, that is, X is sub-Gaussian with $\sigma = (b-a)/2$.*

In the proof we will use the following simple result.

Lemma 5.2.3 (Popoviciu's inequality). *If $X \in [a, b]$ then $\text{var}(X) \leq \left(\frac{b-a}{2}\right)^2$.*

Proof. For every $u \in \mathbb{R}$, $\mathbb{E}(X-u)^2 = \text{var}(X) + (u-\mu)^2$ and so $\text{var}(X) \leq \mathbb{E}(X-u)^2$. Take $u = \frac{a+b}{2}$ then

$$\mathbb{E}(X-u)^2 = \frac{1}{4}\mathbb{E}(X-a+X-b)^2 \leq \frac{1}{4}\mathbb{E}(X-a+b-X)^2 = \left(\frac{b-a}{2}\right)^2.$$

□

Proof of Proposition 5.2.2. Suppose f is the density of X with respect to the underlying measure μ . Consider the cumulant generating function $K(\lambda) = \log \mathbb{E}e^{\lambda X}$ and note that $K(0) = 0$, $K'(0) = \mathbb{E}X$, and $K''(0) = \text{var}(X)$. Directly by definition of $K(\lambda)$

$$f(x; \lambda) = f(x)e^{\lambda x - K(\lambda)}$$

defines a valid density function. Direct calculations (or Proposition 1.3.4) show that the second derivative of $K(\lambda)$ is the variance of the distribution with density $e^{\lambda x - K(\lambda)}f(x)$ ⁵. In particular, $K''(\lambda) = \mathbb{E}_\lambda[X^2] - (\mathbb{E}_\lambda[X])^2$, where $\mathbb{E}_\lambda[g(X)] = \mathbb{E}[g(X)\frac{e^{\lambda X}}{\mathbb{E}[e^{\lambda X}]}] = \int_{\mathbb{R}} g(x)e^{\lambda x - K(\lambda)}f(x)dx$. By Lemma 5.2.3,

$$0 \leq K''(\lambda) = \mathbb{E}_\lambda[X^2] - (\mathbb{E}_\lambda[X])^2 \leq \left(\frac{b-a}{2}\right)^2.$$

By the Taylor's theorem, for every $\lambda \in \mathbb{R}$

$$K(\lambda) = K(0) + K'(0)\lambda + K''(\theta)\frac{\lambda^2}{2}$$

for some θ between 0 and λ . Using what we know, we conclude that

$$K(\lambda) - \lambda\mathbb{E}X \leq \left(\frac{b-a}{2}\right)^2 \frac{\lambda^2}{2} \quad \text{for all } \lambda \in \mathbb{R},$$

which is equivalent to the claimed inequality. □

The following elementary result will be useful.

Lemma 5.2.4. *Suppose that $\mathbf{X} = (X_1, \dots, X_n)$ is a random vector such that X_i is σ_i -sub-Gaussian for $i = 1, \dots, n$. If X_i are independent and $\mathbf{u} \in \mathbb{R}^n$, then the random variable $\mathbf{u}^\top \mathbf{X}$ is sub-Gaussian with parameter $\sigma = \sqrt{\sigma_1^2 u_1^2 + \dots + \sigma_n^2 u_n^2}$.*

⁵ Note that λ is a parameter describing an exponential family where f corresponds to $\lambda = 0$.

Proof. Denote by $\boldsymbol{\mu}$ the mean vector of $\mathbb{E}\mathbf{X}$. Equivalently we must show $K_{\mathbf{u}^\top \mathbf{X}}(\lambda) - \lambda \mathbf{u}^\top \boldsymbol{\mu} \leq \sigma^2 \lambda^2 / 2$. We have

$$K_{\mathbf{u}^\top \mathbf{X}}(\lambda) - \lambda \mathbf{u}^\top \boldsymbol{\mu} = \sum_{i=1}^n (K_{X_i}(\lambda u_i) - \lambda u_i \mu_i) \leq \sum_{i=1}^n \frac{\sigma_i^2 \lambda^2 u_i^2}{2} = \left(\sum_{i=1}^n \sigma_i^2 u_i^2 \right) \frac{\lambda^2}{2},$$

which completes the proof. \square

We list a couple of useful observations that follow directly from Lemma 5.2.4. They show that the σ parameter behaves similarly to the standard deviation.

Corollary 5.2.5. *For a collection of independent σ -sub-Gaussian random variables X_1, \dots, X_n , their average \bar{X}_n is sub-Gaussian with parameter σ / \sqrt{n} .*

Corollary 5.2.6. *If X is σ -sub-Gaussian, then $aX + b$ ($a \in \mathbb{R}, b \in \mathbb{R}$) is sub-Gaussian with parameter $|a|\sigma$. In particular, $-X$ is σ -sub-Gaussian.*

Corollary 5.2.7. *If $\mathbf{X} = (X_1, \dots, X_n)$ is a vector of independent random variables such that each X_i is σ -sub-Gaussian, and $\|\mathbf{u}\| = 1$, then $\mathbf{u}^\top \mathbf{X}$ is σ -sub-Gaussian.*

Suppose that X is σ -sub-Gaussian. By Corollary 5.2.6, $-X$ is also σ -sub-Gaussian. Thus, on the top of (5.5) we also have

$$\mathbb{P}(X - \mu \leq -t) \leq e^{-t^2/2\sigma^2} \quad \text{for all } t \geq 0.$$

Using the union bound, we conclude the following result.

Corollary 5.2.8. *For a σ -sub-Gaussian variable*

$$\mathbb{P}(|X - \mu| \geq t) \leq 2e^{-t^2/2\sigma^2} \quad \text{for all } t \geq 0.$$

The following result applies Chernoff bounds to a sequence of sub-Gaussian variables.

Proposition 5.2.9 (Hoeffding inequality). *If X_i are independent σ_i -sub-Gaussian then for all $t \geq 0$ we have*

$$\mathbb{P}(|\bar{X}_n - \mathbb{E}(\bar{X}_n)| \geq t) \leq 2e^{-t^2 n^2 / (2 \sum_i \sigma_i^2)}$$

Proof. By Lemma 5.2.4, \bar{X}_n is sub-Gaussian with parameter

$$\sigma = \frac{1}{n} \sqrt{\sum_{i=1}^n \sigma_i^2}.$$

Thus, the result follows by Corollary 5.2.8. \square

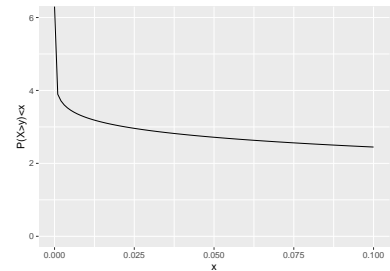
Remark 5.2.10. *If $\sigma_i = \sigma$ for all i we also note that the Hoeffding inequality implies*

$$\mathbb{P}(|\sqrt{n}(\bar{X}_n - \mathbb{E}(\bar{X}_n))| \geq t) \leq 2e^{-t^2/(2\sigma^2)}.$$

In other words, for $\delta \in (0, 1)$,

$$\mathbb{P}\left(|X - \mu| \geq \sigma \sqrt{2 \log\left(\frac{2}{\delta}\right)}\right) \leq \delta.$$

The graph below shows the price we are paying by decreasing $\delta \in (0, 0.1)$ (the x-axis) with respect to the guaranteed bound (y-axis); $\sigma = 1$.



Example 5.2.11 (Sub-Gaussian sequence model). Consider independent Y_1, \dots, Y_n that are σ -sub-Gaussian with means μ_i , where σ^2 is known. Let $S \subset \{1, \dots, n\}$ be the support of the vector $\mu = (\mu_1, \dots, \mu_n)$, that is, $\mu_i \neq 0$ for $i \in S$ and $\mu_i = 0$ for $i \notin S$. One possible way of finding the zero μ_i 's is by considering the LASSO estimator

$$\hat{\mu} = \arg \min_{\mu} \left\{ \frac{1}{2} \|\mathbf{Y} - \mu\|^2 + \lambda \sum_{i=1}^n |\mu_i| \right\}$$

for a fixed $\lambda \geq 0$. It is a simple exercise⁶ to show that

$$\hat{\mu}_i = \begin{cases} Y_i - \lambda & \text{if } Y_i > \lambda, \\ 0 & \text{if } |Y_i| \leq \lambda, \\ Y_i + \lambda & \text{if } Y_i < -\lambda. \end{cases}$$

To analyse this procedure we consider two types of errors

1. Type I: $|Y_i| > \lambda$ for $i \notin S$, and
2. Type II: $|Y_i| \leq \lambda$ for $i \in S$.

We would like to set the threshold λ to control the Family-wise Error Rate. In other words, $Z := \max_{i \notin S} |Y_i| \leq \lambda$ with high probability. Let s be the number of the nonzeros entries in μ then $m_0 := n - s$ is the number of zeros. By the union bound and using the fact that each Y_i for $i \notin S$ is mean zero σ -sub-Gaussian

$$\mathbb{P}(Z \geq \lambda) \leq \sum_{i \notin S} \mathbb{P}(|Y_i| \geq \lambda) \leq 2m_0 e^{-\lambda^2/2\sigma^2}.$$

Suppose that $\lambda = \sigma \sqrt{2 \log(\frac{2m_0}{\alpha})}$ for some $\alpha \in (0, 1)$ then $\mathbb{P}(Z \geq \lambda) \leq \alpha$, which binds the family-wise error rate. Moreover, if $\lambda_n = 2\sigma \sqrt{\log(m_0)}$, then

$$\mathbb{P}(Z \geq \lambda_n) \leq \frac{2}{m_0},$$

which goes to zero as m_0 grows to infinity (with n)⁷. Of course, in practice, we also want λ to be as small as possible to reduce the type II error. We will not discuss this issue here⁸.

Another general type of bound on the moment generating function is the following.

Definition 5.2.12. A random variable with mean $\mu = \mathbb{E}X$ is sub-exponential if there are non-negative parameters ν, α such that

$$\mathbb{E}(e^{\lambda(X-\mu)}) \leq e^{\nu^2 \lambda^2 / 2} \quad \text{for all } |\lambda| < \frac{1}{\alpha}.$$

This condition is rather mild and it is essentially equivalent with the existence of the cumulant generating function in a neighbourhood of zero.

⁶ Check this!

In fact, our analysis applies to any procedure based on thresholding.

⁷ We could compare this procedure with the Bonferroni and the Holm procedures discussed earlier. Note that the latter two have no theoretical guarantees if the underlying distribution is not Gaussian.

⁸ Exercise: Analyze the type II error as a function of $\mu^* = \min_{i \in S} |\mu_i|$

Example 5.2.13. Let $Z \sim N(0,1)$ and let $X = Z^2$. For $\lambda < 1/2$, we have

$$\mathbb{E}(e^{\lambda(X-1)}) = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{\lambda(z^2-1)} e^{-z^2/2} dz = \begin{cases} \frac{e^{-\lambda}}{\sqrt{1-2\lambda}} & \text{if } \lambda < \frac{1}{2}, \\ +\infty & \text{if } \lambda \geq \frac{1}{2}. \end{cases}$$

In particular, since the moment generating function is not defined everywhere, X is not sub-Gaussian. With a bit of calculus we see however that

$$\frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \leq e^{4\lambda^2/2} \quad \text{for all } |\lambda| < \frac{1}{4}$$

and so X is sub-exponential with parameters $(\nu, \alpha) = (2, 4)$.

With essentially the same proof, Lemma 5.2.4 generalizes to sub-exponential variables.

Lemma 5.2.14. Suppose that $\mathbf{X} = (X_1, \dots, X_n)$ is a random vector such that X_i is (ν_i, α_i) -sub-exponential for $i = 1, \dots, n$. If X_i are independent and $\mathbf{u} \in \mathbb{R}^n$, then the random variable $\mathbf{u}^\top \mathbf{X}$ is sub-exponential with parameters $\nu = \sqrt{\nu_1^2 u_1^2 + \dots + \nu_n^2 u_n^2}$ and $\alpha = \max_i |u_i| \alpha_i$.

We obtain simple concentration inequalities for sub-exponential variables.

Proposition 5.2.15. Suppose that X is sub-exponential with parameters (ν, α) and $t \geq 0$. Then

$$\mathbb{P}(X - \mu \geq t) \leq \begin{cases} e^{-\frac{t^2}{2\nu^2}} & \text{if } 0 \leq t \leq \frac{\nu^2}{\alpha}, \\ e^{-\frac{t}{2\alpha}} & \text{if } t > \frac{\nu^2}{\alpha}. \end{cases}$$

and

$$\mathbb{P}(|X - \mu| \geq t) \leq \begin{cases} 2e^{-\frac{t^2}{2\nu^2}} & \text{if } 0 \leq t \leq \frac{\nu^2}{\alpha}, \\ 2e^{-\frac{t}{2\alpha}} & \text{if } t > \frac{\nu^2}{\alpha}. \end{cases}$$

Proof. Using the Chernoff bound we know that

$$\mathbb{P}(X - \mu \geq t) \leq \inf_{\lambda \in [0, \frac{1}{\alpha})} \{\mathbb{E}e^{\lambda(X - \mathbb{E}X)} e^{-\lambda t}\} \leq \inf_{\lambda \in [0, \frac{1}{\alpha})} \{e^{\nu^2 \lambda^2 / 2 - \lambda t}\}.$$

The global optimum of $e^{\nu^2 \lambda^2 / 2 - \lambda t}$ is $\lambda^* = t / \nu^2 \geq 0$. Thus, if $t / \nu^2 < 1 / \alpha$ (equiv. $\nu^2 > t\alpha$), we get the sub-Gaussian bound $e^{-t^2 / (2\nu^2)}$.

Otherwise, if $\nu^2 \leq t\alpha$, the infimum is obtained on the boundary $\lambda^* = 1 / \alpha$ and it is equal to $e^{\nu^2 / (2\alpha^2) - t / \alpha}$. Note however that the fact that $\nu^2 \leq t\alpha$, allows to bound this by $e^{-t / (2\alpha)}$. \square

Lemma 5.2.14 and Proposition 5.2.15 give now a handful of useful results. For example, if X_i are independent sub-exponential with parameters (ν, α) then \bar{X}_n is sub-exponential with parameters $(\frac{\nu}{\sqrt{n}}, \frac{\alpha}{n})$ and so

$$\mathbb{P}(|\bar{X}_n - \mu| \geq t) \leq \begin{cases} 2e^{-\frac{t^2 n}{2\nu^2}} & \text{if } 0 \leq t \leq \frac{\nu^2}{\alpha}, \\ 2e^{-\frac{tn}{2\alpha}} & \text{if } t > \frac{\nu^2}{\alpha}. \end{cases}$$

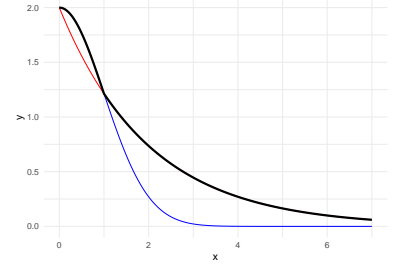


Figure 5.1: Illustration of the sub-exponential bounds for $\nu = \alpha = 1$. The blue line represents $2e^{-t^2/2}$. The red line represents $2e^{-t/2}$. In black their pointwise maximum.

5.3 Martingale-based methods

We will now briefly discuss techniques related to martingale representations. Recall the basic set-up of the martingale theory in Appendix C.2. Let X_1, \dots, X_n be independent random variables and $Z = f(X_1, \dots, X_n)$. Denote

$$\mathbb{E}_i[\cdot] = \mathbb{E}[\cdot | X_1, \dots, X_i].$$

Thus $\mathbb{E}_0 Z = \mathbb{E}Z$ and $\mathbb{E}_n Z = Z$. Moreover, if $i < j$ then $\mathbb{E}_i \mathbb{E}_j Z = \mathbb{E}_j \mathbb{E}_i Z = \mathbb{E}_i Z$. It is easy to see that $Y_k = \mathbb{E}_k Z$ forms a martingale sequence if $\mathbb{E}|Z| < \infty$. Indeed, by Jensen's inequality

$$\mathbb{E}(|\mathbb{E}_k Z|) \leq \mathbb{E}\mathbb{E}_k |Z| = \mathbb{E}|Z| < \infty.$$

This is called the Doob's martingale.

Writing

$$\Delta_i = \mathbb{E}_i Z - \mathbb{E}_{i-1} Z$$

we have

$$Z - \mathbb{E}Z = \sum_{i=1}^n \Delta_i. \quad (5.7)$$

This is the Doob martingale representation of Z . In particular, we write $Z - \mathbb{E}Z = \sum_{i=1}^n \Delta_i$, where Δ_i for the corresponding Doob martingale difference sequence.

Theorem 5.3.1. *Let $\{\Delta_i, \mathcal{F}_i\}$ be a martingale difference sequence as above and suppose that $\mathbb{E}_{k-1}[e^{\lambda \Delta_k}] \leq e^{\lambda^2 v_k^2 / 2}$ almost surely for any $|\lambda| < 1/\alpha_k$ (a form of a sub-exponential condition). Then $Z - \mathbb{E}Z = \sum_{i=1}^n \Delta_i$ is subexponential with parameters $(\nu_*, \alpha_*) = (\|\mathbf{v}\|, \|\boldsymbol{\alpha}\|_\infty)$. In particular,*

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq \begin{cases} 2e^{-\frac{t^2}{2\nu_*^2}} & \text{if } 0 \leq t \leq \frac{\nu_*^2}{\alpha_*}, \\ 2e^{-\frac{t}{2\alpha_*}} & \text{if } t > \frac{\nu_*^2}{\alpha_*}. \end{cases}$$

Proof. The first part can be directly shown by recursive conditioning. Indeed, we first write

$$\mathbb{E}e^{\lambda(Z - \mathbb{E}Z)} = \mathbb{E}e^{\lambda \sum_{i=1}^n \Delta_i} = \mathbb{E}\mathbb{E}_{n-1} e^{\lambda \sum_{i=1}^n \Delta_i} = \mathbb{E}e^{\lambda \sum_{i=1}^{n-1} \Delta_i} \mathbb{E}_{n-1} e^{\lambda \Delta_n}$$

and this is bounded for all $|\lambda| < 1/\alpha_n$ by

$$e^{\lambda^2 \nu_n^2 / 2} \mathbb{E}e^{\lambda \sum_{i=1}^{n-1} \Delta_i}.$$

Proceeding in a similar fashion, we get the conclusion. The second of the theorem part follows by the first part and Proposition 5.2.15. \square

A very useful version of this result is when each Δ_i is bounded, $\Delta_k \in [a_k, b_k]$, in which case it is also sub-Gaussian (also conditionally on \mathcal{F}_{k-1}).

Definition 5.3.2. A function $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfies **bounded differences inequality**, if there exist constants L_1, \dots, L_n such that

$$|f(x) - f(x')| \leq L_k,$$

whenever x, x' differ only in the k -th coordinate.

A canonical example of such a function is

$$f(x) = \sum_{i=1}^n f_i(x_i),$$

where all f_i are bounded functions.

Proposition 5.3.3 (Bounded differences inequality). *Suppose that $Z = f(X_1, \dots, X_n)$, where f satisfies the bounded differences inequality with parameters L_1, \dots, L_n and such that the random vector $X = (X_1, \dots, X_n)$ has independent components. Then Z is sub-Gaussian with parameter*

$$\sigma = \frac{1}{2} \sqrt{L_1^2 + \dots + L_n^2}.$$

In particular, for all $t \geq 0$

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq 2e^{-\frac{t^2}{2\sigma^2}} = 2e^{-\frac{2t^2}{\sum_{i=1}^n L_i^2}}.$$

Proof. Fix i and take a look at Δ_i . Let

$$a_i := \inf_x \mathbb{E}[Z | X_1, \dots, X_{i-1}, x] - \mathbb{E}_{i-1} Z$$

and

$$b_i := \sup_x \mathbb{E}[Z | X_1, \dots, X_{i-1}, x] - \mathbb{E}_{i-1} Z.$$

It is clear that $a_i \leq \Delta_i \leq b_i$ almost surely. To use Proposition 5.2.2 we need to bound $b_i - a_i$. For that, note first that

$$\mathbb{E}[f(X_1, \dots, X_i, \dots, X_n) | X_1, \dots, X_{i-1}, X_i = x] = \mathbb{E}[f(X_1, \dots, x, \dots, X_n) | X_1, \dots, X_{i-1}],$$

which follows by independence of all X_i 's (in both cases the integral is with respect to the marginal distribution of (X_{i+1}, \dots, X_n)).

We have

$$\begin{aligned} b_i - a_i &= \sup_x \mathbb{E}[Z | X_1, \dots, X_{i-1}, x] - \inf_x \mathbb{E}[Z | X_1, \dots, X_{i-1}, x] \\ &\leq \sup_{x,y} \left| \mathbb{E}[Z | X_1, \dots, X_{i-1}, x] - \mathbb{E}[Z | X_1, \dots, X_{i-1}, y] \right| \\ &= \sup_{x,y} \left| \mathbb{E}_{i-1} [f(X_1, \dots, X_{i-1}, x, X_{i+1}, \dots, X_n) - f(X_1, \dots, X_{i-1}, y, X_{i+1}, \dots, X_n)] \right| \\ &\leq L_i. \end{aligned}$$

Since $\Delta_i \in [a_i, b_i]$ and $b_i - a_i \leq L_i$, by Proposition 5.2.2,

$$\mathbb{E}_{i-1} e^{\lambda \Delta_i} \leq e^{\lambda^2 v_i^2 / 2}, \quad \text{for } v_i = L_i / 2 \text{ and all } \lambda \in \mathbb{R}. \quad (5.8)$$

Now we can use Theorem 5.3.1 to conclude that $Z - \mathbb{E}Z$ (and so also Z) is sub-Gaussian with parameter $\frac{1}{2}\sqrt{L_1^2 + \dots + L_n^2}$. The probability bound then follows from (5.5). \square

Remark 5.3.4. Note that there was an alternative route to simply observe that $|f(x) - f(y)| \leq \sum_k L_k$ for all $x, y \in \mathcal{X}^n$. In other words, Z is bounded and so sub-Gaussian with parameter $\sigma = \sum_k L_k$ even without assuming independence of X_1, \dots, X_n . However, $\frac{1}{2}\sqrt{\sum_k L_k^2} < \sum_k L_k$ and this difference can be critical.

Example 5.3.5 (Kernel density estimation). Let X_1, \dots, X_n be i.i.d. real samples drawn according to some density ϕ . The kernel density estimate is

$$\phi_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right),$$

where $h > 0$, and K is a nonnegative “kernel” $\int K = 1$. The L_1 -error is

$$Z = f(X_1, \dots, X_n) = \int_{\mathcal{X}} |\phi_n(x) - \phi(x)| dx.$$

It is easy to see that, if X, X' differ only in the i -th coordinate then

$$|f(X) - f(X')| \leq \frac{1}{nh} \int \left| K\left(\frac{x - X_i}{h}\right) - K\left(\frac{x - X'_i}{h}\right) \right| dx \leq \frac{2}{n}.$$

By Proposition 5.3.3,

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq 2e^{-\frac{t^2 n}{2}} \quad \text{for all } t \geq 0.$$

Note however that to analyze the quality of the kernel density estimator, we need bounds on Z directly. For that we need to separately study $\mathbb{E}Z$.⁹

Example 5.3.6 (Uniform deviations). Let \mathcal{A} be a collection of measurable subsets of \mathcal{X} and let X, X_1, \dots, X_n be random points drawn from \mathcal{X} i.i.d. Let

$$\mathbb{P}(A) = \mathbb{P}\{X \in A\} \quad \text{and} \quad \mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \in A\}.$$

Let $Z = f(X_1, \dots, X_n) = \sup_{A \in \mathcal{A}} |\mathbb{P}(A) - \mathbb{P}_n(A)|$. If X, X' differ only in the i -th coordinate then, for every $A \in \mathcal{A}$,

$$\left| |\mathbb{P}(A) - \mathbb{P}_n(A)| - |\mathbb{P}(A) - \mathbb{P}'_n(A)| \right| \leq |\mathbb{P}_n(A) - \mathbb{P}'_n(A)| = \frac{1}{n} |\mathbb{1}\{X_i \in A\} - \mathbb{1}\{X'_i \in A\}| \leq \frac{1}{n}.$$

We conclude that $|f(X) - f(X')| \leq 1/n$. By Proposition 5.3.3,

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq 2e^{-2t^2 n} \quad \text{for all } t \geq 0.$$

This holds irrespective of the underlying distribution or richness of \mathcal{A} . As in the previous example, a major issue can be to understand better $\mathbb{E}Z$.

⁹ This is well studied in the literature. See for example Section II.4.2 in P. P. B. Eggermont and V. N. LaRiccia. *Maximum penalized likelihood estimation. Vol. I.* Springer Series in Statistics. Springer-Verlag, New York, 2001

Example 5.3.7 (U-statistics). Let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a symmetric function. Given an i.i.d. sequence X_i of random variables, the quantity

$$Z := \frac{1}{\binom{n}{2}} \sum_{i < j} g(X_i, X_j).$$

is known as a pairwise U-statistic. For instance, if $g(s, t) = |s - t|$, then U is an unbiased estimator of the mean absolute pairwise deviation $\mathbb{E}(|X_1 - X_2|)$. Suppose g is bounded, say $\|g\|_\infty \leq b$. Writing $Z =: f(X_1, \dots, X_n)$. For any two $x, x' \in \mathbb{R}^n$ that differ only in the i -th coordinate we have

$$|f(x) - f(x')| \leq \frac{1}{\binom{n}{2}} \sum_{j \neq i} |g(x_i, x_j) - g(x'_i, x_j)| \leq \frac{(n-1)2b}{\binom{n}{2}} = \frac{4b}{n}.$$

By Proposition 5.3.3,

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq 2e^{-\frac{t^2 n}{8b^2}} \quad \text{for all } t \geq 0.$$

In particular, Z is a consistent estimator of $\mathbb{E}Z = \mathbb{E}g(X_1, X_2)$.

The following example will become important later in the lecture. This is also our first example for which the constants L_1, \dots, L_n may be different.

Example 5.3.8 (Rademacher complexity). Let $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ be a vector of independent Rademacher random variables (values ± 1 with equal probability). Given a collection of vectors $\mathcal{A} \subseteq \mathbb{R}^n$, define the random variable

$$Z(\mathcal{A}) := \sup_{a \in \mathcal{A}} \langle a, \varepsilon \rangle.$$

Although we do not use it here, note that $Z(\mathcal{A})$ is a convex function of ε .

The random variable $Z = Z(\mathcal{A})$ measures the size of \mathcal{A} in a certain sense, and its expectation $\mathcal{R}(\mathcal{A})$ is known as the **Rademacher complexity** of the set \mathcal{A} . Suppose \mathcal{A} is bounded. We will use Proposition 5.3.3 to show that $Z = f(\varepsilon_1, \dots, \varepsilon_n)$ is sub-Gaussian. Note that for every $a \in \mathcal{A}$, if $\varepsilon, \varepsilon'$ differ only in the i -th coordinate, then

$$\langle a, \varepsilon \rangle - \sup_{b \in \mathcal{A}} \langle b, \varepsilon' \rangle = \inf_{b \in \mathcal{A}} (\langle a, \varepsilon \rangle - \langle b, \varepsilon' \rangle) \leq \langle a, \varepsilon \rangle - \langle a, \varepsilon' \rangle = a_i(\varepsilon_i - \varepsilon'_i) \leq 2|a_i|.$$

Thus, we can take the supremum over $a \in \mathcal{A}$ to conclude that

$$f(\varepsilon) - f(\varepsilon') = \sup_{a \in \mathcal{A}} \langle a, \varepsilon \rangle - \sup_{b \in \mathcal{A}} \langle b, \varepsilon' \rangle \leq 2 \sup_{a \in \mathcal{A}} |a_i| =: L_i.$$

The argument is symmetric if we swap ε and ε' and so f satisfies the bounded differences inequality. By Proposition 5.3.3, $Z(\mathcal{A})$ is then sub-Gaussian with parameter $2\sqrt{\sum_i \sup_{a \in \mathcal{A}} a_i^2}$. We remark that, using different techniques, this sub-Gaussianity parameter can be reduced to $\sup_{a \in \mathcal{A}} \|a\|$.

An important application of this theory is for the variable $Z = \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$, as defined in (4.1), at least in the case when \mathcal{F} is uniformly bounded.

Theorem 5.3.9. *Assume \mathcal{F} is uniformly b -bounded, that is, $\|f\|_{\infty} \leq b$ for all $f \in \mathcal{F}$. We have*

$$\mathbb{P}(Z \geq \mathbb{E}Z + t) \leq \exp\left\{-\frac{nt^2}{2b^2}\right\}.$$

Proof. The function $g(x_1, \dots, x_n) = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}f(X) \right|$ satisfies the bounded difference property with $L_i = \frac{2b}{n}$. \square

5.4 Lipschitz functions of Gaussian variables

Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is **Lipschitz with parameter L** if

$$|f(x) - f(y)| \leq L\|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n.$$

Theorem 5.4.1. *Let $X = (X_1, \dots, X_n) \sim N_n(0, I_n)$ and let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be Lipschitz with parameter L . Then the variable $f(X)$ is L -sub-Gaussian, and hence*

$$\mathbb{P}(|f(X) - \mathbb{E}f(X)| \geq t) \leq 2e^{-\frac{t^2}{2L^2}} \quad \text{for all } t \geq 0.$$

We will prove this result with a slightly worse constant and assuming differentiability (Lipschitz functions are differentiable almost everywhere though). In this case $\|\nabla f(x)\| \leq L$ for all x (easy exercise, e.g. consider the directional derivative of $f(x)$ in the direction $\nabla f(x)$).

Proof. We have $f(X) - \mathbb{E}f(X) = f(X) - \mathbb{E}_{X'}f(X')$, where X' is an independent copy of X . By the Jensen's inequality

$$\mathbb{E}_X e^{\lambda(f(X) - \mathbb{E}_{X'}f(X'))} \leq \mathbb{E}_{X, X'} e^{\lambda(f(X) - f(X'))}. \quad (5.9)$$

Suppose f is differentiable then, by the fundamental theorem of calculus,

$$\begin{aligned} f(X) - f(X') &= \int_0^{\pi/2} \frac{d}{d\theta} f(X \sin \theta + X' \cos \theta) d\theta \\ &= \int_0^{\pi/2} \langle \nabla f(X \sin \theta + X' \cos \theta), X \cos \theta - X' \sin \theta \rangle d\theta \end{aligned}$$

Note that the variables $X_{\theta} := X \sin \theta + X' \cos \theta$ and $X'_{\theta} := X \cos \theta - X' \sin \theta$ are independent standard normal and so the distribution of (X, X') is the same as the distribution of $(X_{\theta}, X'_{\theta})$. The right-hand side in (5.9) can be rewritten as

$$\mathbb{E}_{X, X'} e^{\lambda \int_0^{\pi/2} \langle \nabla f(X_{\theta}), X'_{\theta} \rangle d\theta}.$$

Note that the sub-Gaussian parameter is completely dimension-free! For this reason this result is sometimes referred to as "Dimension Free Concentration Inequality".

If θ is a random variable uniformly distributed on $(0, \frac{\pi}{2})$ then

$$\int_0^{\pi/2} \langle \nabla f(X_\theta), X'_\theta \rangle d\theta = \frac{\pi}{2} \mathbb{E}_\theta \langle \nabla f(X_\theta), X'_\theta \rangle.$$

Using the Jensen's inequality again (and Fubini's theorem), we get

$$\mathbb{E}_{X, X'} e^{\lambda \int_0^{\pi/2} \langle \nabla f(X_\theta), X'_\theta \rangle d\theta} = \mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \mathbb{E}_\theta \langle \nabla f(X_\theta), X'_\theta \rangle} \leq \mathbb{E}_\theta \mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \langle \nabla f(X_\theta), X'_\theta \rangle}.$$

For any fixed θ , (X_θ, X'_θ) has the same distribution as (X, X') . It then follows that

$$\mathbb{E}_\theta \mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \langle \nabla f(X_\theta), X'_\theta \rangle} = \mathbb{E}_\theta \mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \langle \nabla f(X), X' \rangle} = \mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \langle \nabla f(X), X' \rangle}.$$

Note that for fixed X , the variable $\langle \nabla f(X), X' \rangle$ is Gaussian with mean zero and variance $\|\nabla f(X)\|^2$ and hence

$$\mathbb{E}_{X, X'} e^{\lambda \frac{\pi}{2} \langle \nabla f(X), X' \rangle} \leq \mathbb{E}_X e^{\lambda^2 \pi^2 \|\nabla f(X)\|^2 / 8} \leq e^{\lambda^2 \pi^2 L^2 / 8}.$$

This calculation shows that

$$\mathbb{E}_X e^{\lambda(f(X) - \mathbb{E}_{X'} f(X'))} \leq e^{\lambda^2 \pi^2 L^2 / 8}$$

or, in other words, that $f(X) - \mathbb{E}f(X)$ is sub-Gaussian with parameter $\sigma = \frac{\pi L}{2}$, which is slightly more than the claimed L . \square

This result is useful for a broad range of problems.

Example 5.4.2 (χ^2 concentration). For a given vector $Z = (Z_1, \dots, Z_n)$ of i.i.d. standard normal variables, we have $Y := \sum_{i=1}^n Z_i^2 = \|Z\|^2 \sim \chi_n^2$. The most direct way to obtain tail bounds on Y was given in Example 5.1.4. Alternatively we can use Theorem 5.4.1. Define $V = \sqrt{Y/n} = \|Z\|/\sqrt{n}$. Since the Euclidean norm is a 1-Lipschitz function, we get

$$\mathbb{P}(V - \mathbb{E}V \geq t) \leq e^{-nt^2/2}.$$

To obtain bounds on $Y - \mathbb{E}Y$ note that, by Jensen's inequality,

$$\mathbb{E}[V] \leq \sqrt{\mathbb{E}[V^2]} = 1.$$

Thus

$$\mathbb{P}(V - \mathbb{E}V \geq t) \geq \mathbb{P}(V - 1 \geq t) = \mathbb{P}(Y/n \geq (t+1)^2).$$

Putting everything together we conclude

$$\mathbb{P}(Y/n \geq (1+t)^2) \leq e^{-nt^2/2} \quad \text{for all } t \geq 0.$$

Example 5.4.3 (Order statistics). Given a random vector (X_1, \dots, X_n) , its order statistics are obtained by reordering its entries in a non-decreasing manner as

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n-1)} \leq X_{(n)}.$$

As particular cases, we have $X_{(n)} = \max_i X_i$ and $X_{(1)} = \min_i X_i$. Given another random vector (Y_1, \dots, Y_n) , it can be shown that $|X_{(k)} - Y_{(k)}| \leq \|X - Y\|$ for all $k = 1, \dots, n$ ¹⁰, so that each order statistic is a 1-Lipschitz function. Consequently, when X is a standard Gaussian random vector, Theorem 5.4.1 implies that

$$\mathbb{P}(|X_{(i)} - \mathbb{E}X_{(i)}| \geq t) \leq 2e^{-t^2/2} \quad \text{for all } t \geq 0.$$

Example 5.4.4 (Singular values of Gaussian random matrices). For $n \geq d$, consider the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, and denote by

$$\sigma_1(\mathbf{X}) \geq \sigma_2(\mathbf{X}) \geq \dots \geq \sigma_d(\mathbf{X}) \geq 0$$

the singular values of \mathbf{X} . By Weyl's theorem

$$\max_{i=1, \dots, d} |\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{Y})| \leq \|\mathbf{X} - \mathbf{Y}\| \leq \|\mathbf{X} - \mathbf{Y}\|_F.$$

In other words, each singular value $\sigma_i(\mathbf{X})$ is a 1-Lipschitz function of \mathbf{X} . Suppose now that \mathbf{W} is random with independent standard normal entries. Theorem 5.4.1 implies that

$$\mathbb{P}(|\sigma_i(\mathbf{W}) - \mathbb{E}\sigma_i(\mathbf{W})| \geq t) \leq 2e^{-t^2/2} \quad \text{for all } t \geq 0,$$

or in other words

$$\mathbb{P}(|\sigma_i(\frac{1}{\sqrt{n}}\mathbf{W}) - \mathbb{E}\sigma_i(\frac{1}{\sqrt{n}}\mathbf{W})| \geq t) \leq 2e^{-t^2n/2} \quad \text{for all } t \geq 0.$$

If \mathbf{X} has i.i.d. rows from $N_d(0, \Sigma)$ then $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$ and the sample covariance satisfies

$$\hat{\Sigma} = \frac{1}{n}\mathbf{X}^\top\mathbf{X} = \sqrt{\Sigma}(\frac{1}{\sqrt{n}}\mathbf{W})^\top(\frac{1}{\sqrt{n}}\mathbf{W})\sqrt{\Sigma}.$$

We can exploit this to obtain some bounds on $\|\hat{\Sigma} - \Sigma\|$ (see Chapter 6 in ¹¹).

Example 5.4.5 (Gaussian width). Let $W = (W_1, \dots, W_n)$ be the n -dimensional standard Gaussian vector. Given a collection of vectors $\mathcal{A} \subset \mathbb{R}^n$, define the random variable

$$Z = Z(\mathcal{A}) := \sup_{a \in \mathcal{A}} \langle a, W \rangle.$$

The variable Z is one way of measuring the size of the set \mathcal{A} , which is called the Gaussian width. We view Z as a function $(w_1, \dots, w_n) \mapsto f(w_1, \dots, w_n)$. Fixing a , we get

$$\langle a, w \rangle - \langle a, w' \rangle \leq \|a\| \|w - w'\| \leq \sup_{a \in \mathcal{A}} \|a\| \|w - w'\|.$$

Then

$$\sup_{a \in \mathcal{A}} \langle a, w \rangle - \sup_{a \in \mathcal{A}} \langle a, w' \rangle \leq \sup_{a \in \mathcal{A}} \|a\| \|w - w'\|,$$

¹⁰ The fact that $|X_k - Y_k| \leq \|X - Y\|$ for every k is trivial. But this statement is more subtle!

¹¹ Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019

and similarly

$$\sup_{a \in \mathcal{A}} \langle a, w' \rangle - \sup_{a \in \mathcal{A}} \langle a, w \rangle \leq \sup_{a \in \mathcal{A}} \|a\| \|w - w'\|.$$

Let $D(\mathcal{A}) := \sup_{a \in \mathcal{A}} \|a\|$. Thus $f(w)$ is $D(\mathcal{A})$ –Lipschitz and this

$$\mathbb{P}(|Z - \mathbb{E}Z| \geq t) \leq 2e^{-\frac{t^2}{2D^2(\mathcal{A})}}.$$

5.5 Exercises

Exercise 5.5.1. Provide a non-negative random variable for which the Markov inequality holds as equality (for a fixed t). Is there a random variable for which it holds as equality for every $t \geq 0$?

Exercise 5.5.2. Show that for $U \sim N(0, 1)$ it holds that

$$\mathbb{P}(U \geq t) \leq \frac{1}{\sqrt{2\pi}} \frac{1}{t} e^{-t^2/2}.$$

Hint: If $\phi(u)$ is the density of U , show $\phi'(u) + u\phi(u) = 0$ for all $u \in \mathbb{R}$. Conclude that $\int_t^\infty \phi(u) du \leq \frac{1}{t}\phi(t)$.

Exercise 5.5.3. Show that in Proposition 5.1.3 the same way we get

$$\log \mathbb{P}(Z - \mathbb{E}Z \leq -t) \leq \inf_{\lambda \in (-b, 0]} \{\log \mathbb{E} e^{\lambda(Z - \mathbb{E}Z)} + \lambda t\}.$$

Conclude bounds on $\mathbb{P}(|Z - \mathbb{E}Z| \geq t)$.

Exercise 5.5.4. Use Exercise 5.5.3 to show that in the Example 5.1.4

$$\mathbb{P}(\bar{Y}_n - 1 \leq -t) \leq (1 - t)^{n/2} e^{nt/2}.$$

Conclude for $t \in [0, 1]$ that

$$\mathbb{P}(|\bar{Y}_n - 1| \geq t) \leq 2e^{-nt^2/8}. \quad (5.10)$$

Exercise 5.5.5. Use Exercise 5.5.3 to show that for all $t \geq 0$

$$\mathbb{P}(X - \mu \leq -t) \leq e^{-t^2/2\sigma^2} \quad \text{and} \quad \mathbb{P}(|X - \mu| \geq t) \leq 2e^{-t^2/2\sigma^2}.$$

Exercise 5.5.6. Suppose X is σ -sub-Gaussian with mean μ . Show that $\text{var}(X) \leq \sigma^2$.

Exercise 5.5.7. Suppose $X_i, i = 1, \dots, N$ are zero-mean σ -sub-gaussian random variables (not necessarily independent). Use Proposition 6.1.1 to show that

$$\mathbb{P}(\max_i X_i - \mathbb{E} \max_i X_i \geq t) \leq e^{-t^2/\sigma^2} \quad \text{for } t < 2\sigma\sqrt{2 \log N}.$$

6

More advanced techniques (1-2 weeks)

6.1 Maximal inequalities

Lemma 5.2.4 and Lemma 5.2.14 provide concentration bounds on linear combinations of sub-Gaussian or sub-exponential variables. In many instances, we will be interested in controlling the maximum over the parameters of such linear combinations. The main motivation is in empirical risk minimization (see (2.1)) but many other applications exist. The purpose of this section is to present such results.

We begin by the simplest case possible: the maximum over a finite set.

Proposition 6.1.1. *Suppose X_1, \dots, X_N are zero-mean σ -sub-Gaussian then*

$$\mathbb{E} \max_{i=1, \dots, N} X_i \leq \sigma \sqrt{2 \log N} \quad \text{and} \quad \mathbb{E} \max_{i=1, \dots, N} |X_i| \leq \sigma \sqrt{2 \log(2N)}.$$

Moreover, for any $t > 0$

$$\begin{aligned} \mathbb{P}(\max_{i=1, \dots, N} X_i \geq t) &\leq N e^{-t^2/(2\sigma^2)} \\ \mathbb{P}(\max_{i=1, \dots, N} |X_i| \geq t) &\leq 2N e^{-t^2/(2\sigma^2)}. \end{aligned} \tag{6.1}$$

In other words, for $\delta \in (0, 1)$,

$$\mathbb{P}\left(\max_i X_i \geq \sigma \sqrt{2 \log\left(\frac{N}{\delta}\right)}\right) \leq \delta.$$

Proof. Take $\lambda > 0$. Then we get

$$\begin{aligned} e^{\lambda \mathbb{E} \max_i X_i} &\stackrel{\text{Jensen}}{\leq} \mathbb{E} e^{\lambda \max_i X_i} = \mathbb{E} \max_i e^{\lambda X_i} \\ &\leq \sum_{i=1}^N \mathbb{E} e^{\lambda X_i} \leq N e^{\sigma^2 \lambda^2 / 2}. \end{aligned}$$

By taking logs and optimizing with respect to λ we get that the optimal $\lambda = \sqrt{2 \log N} / \sigma$. Plugging this back, we obtain

$$\mathbb{E} \max_i X_i \leq \frac{\log N}{\lambda} + \frac{\sigma^2 \lambda}{2} \leq \sigma \sqrt{2 \log N}.$$

The analogous bound for the absolute values follows from the fact that $\max_i |X_i| = \max\{X_1, -X_1, \dots, X_N, -X_N\}$. Finally, the probability bounds follow directly by the union bound and σ -sub-Gaussianity of X_i . \square

Figure 6.1 offers an illustration of Proposition 6.1.1. The blue curve depicts a Monte Carlo estimate of the maximum of N standard normal variables ($\sigma = 1$) and its theoretical bound $\sqrt{2 \log(N)}$ is plotted in red. The bound reveals the right rate. Note also that in Proposition 6.1.1 we did not require that X_i are independent!

Let $\mathbf{X} = (X_1, \dots, X_N)$ be a vector of independent σ -sub-Gaussian variables. It is often important to analyze suprema of random functions over a general set $A \subseteq \mathbb{R}^N$:

$$\sup_{\mathbf{u} \in A} \mathbf{u}^\top \mathbf{X} \quad (6.2)$$

By Corollary 5.2.7, each $\mathbf{u}^\top \mathbf{X}$ is $\sigma \|\mathbf{u}\|$ -sub-Gaussian. If A is finite, we can use Proposition 6.1.1 to get

$$\mathbb{E} \max_{\mathbf{u} \in A} \mathbf{u}^\top \mathbf{X} \leq \sigma \operatorname{diam}(A) \sqrt{2 \log |A|},$$

where $\operatorname{diam}(A) = \max_{\mathbf{u} \in A} \|\mathbf{u}\|$.

Sometimes, even if A is not finite, there is a trivial reduction to the finite case. Suppose that $P \subseteq \mathbb{R}^d$ is a polytope, that is, a set of convex combinations of some fixed points v_1, \dots, v_N . In this case we have

$$\sup_{\mathbf{u} \in P} \mathbf{u}^\top \mathbf{X} = \max_{i=1, \dots, N} v_i^\top \mathbf{X}$$

and now we can use Proposition 6.1.1 to obtain bounds on $\mathbb{E} \sup_{\mathbf{u} \in P} \mathbf{u}^\top \mathbf{X}$ and $\mathbb{P}(\sup_{\mathbf{u} \in P} \mathbf{u}^\top \mathbf{X} \geq t)$. Of particular interest are polytopes that have a small number of vertices. A primary example is the ℓ_1 ball of \mathbb{R}^d defined by

$$\mathbb{B}_1 = \{x \in \mathbb{R}^d : \sum_{i=1}^d |x_i| \leq 1\},$$

which has exactly $2d$ vertices.

If A is bounded, we still have a principled way to study the supremum in (6.2). To simplify the discussion we first present the main ideas in the specific case of the ℓ_2 ball

$$\mathbb{B}_2 = \{x \in \mathbb{R}^d : \sum_{i=1}^d x_i^2 \leq 1\}.$$

The general idea is to cover \mathbb{B}_2 with a finite set of points such that the maximum over this finite set is of the same order as the maximum over the entire ball.

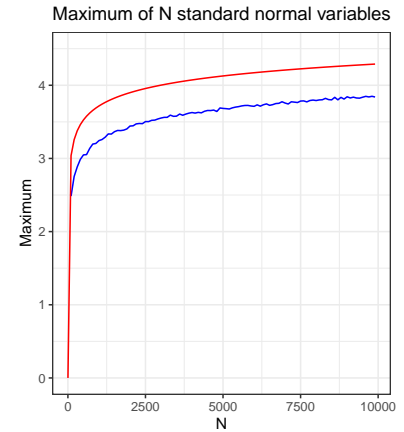


Figure 6.1: Illustration of the bound in Proposition 6.1.1.

← Exercise 5.5.7

Definition 6.1.2. Fix $K \subset \mathbb{R}^d$ and $\epsilon > 0$. A set \mathcal{N} is called an ϵ -cover of K , if $\mathcal{N} \subseteq K$ and for any $z \in K$, there exists $x \in \mathcal{N}$ such that $\|x - z\| \leq \epsilon$.

Definition 6.1.3. The ϵ -covering number of K is

$$N(\epsilon, K) := \inf\{N \in \mathbb{N} : \exists \text{ an } \epsilon\text{-cover of } K \text{ of size } N\}$$

The following lemma gives an upper bound on the size of the smallest ϵ -cover of \mathbb{B}_2 .

Lemma 6.1.4. Fix $\epsilon \in (0, 1)$. Then the unit Euclidean ball \mathbb{B}_2 has an ϵ -cover \mathcal{N} of cardinality $|\mathcal{N}| \leq (3/\epsilon)^d$.

Proof. Consider the following iterative construction of the ϵ -cover. Choose $x_1 = 0$. For any $i \geq 2$, take x_i to be any point in \mathbb{B}_2 such that $\|x_i - x_j\| > \epsilon$ for all $j < i$. If no such point exists, stop the procedure¹. This will create an ϵ -cover. We now control its size.

Since $\|x - y\| > \epsilon$ for all $x, y \in \mathcal{N}$, the Euclidean balls of radius $\epsilon/2$ and centered at points of \mathcal{N} are disjoint. Moreover,

$$\bigcup_{z \in \mathcal{N}} \{z + \frac{\epsilon}{2}\mathbb{B}_2\} \subset (1 + \frac{\epsilon}{2})\mathbb{B}_2.$$

Thus, measuring volumes, we get

$$\text{vol}((1 + \frac{\epsilon}{2})\mathbb{B}_2) \geq \text{vol}(\bigcup_{z \in \mathcal{N}} \{z + \frac{\epsilon}{2}\mathbb{B}_2\}) = |\mathcal{N}| \text{vol}(\frac{\epsilon}{2}\mathbb{B}_2)$$

or equivalently

$$(1 + \frac{\epsilon}{2})^d \geq |\mathcal{N}|(\frac{\epsilon}{2})^d,$$

which gives the following bound

$$|\mathcal{N}| \leq (1 + \frac{2}{\epsilon})^d \leq (\frac{3}{\epsilon})^d$$

□

Recall from Corollary 5.2.7 that if $\mathbf{X} = (X_1, \dots, X_n)$ is a vector of independent σ -sub-Gaussian random variables then for any unit vector \mathbf{u} the variable $\mathbf{u}^\top \mathbf{X}$ is σ -sub-Gaussian. We can then introduce the following general definition of sub-Gaussian vectors.

Definition 6.1.5. A random vector \mathbf{X} is called σ -sub-Gaussian if $\mathbf{u}^\top \mathbf{X}$ is σ -sub-Gaussian for any unit vector \mathbf{u} .

Theorem 6.1.6. Let $\mathbf{X} \in \mathbb{R}^d$ be a mean-zero σ -sub-Gaussian vector. Then

$$\mathbb{E}\|\mathbf{X}\| = \mathbb{E}(\max_{\mathbf{u} \in \mathbb{B}_2} \mathbf{u}^\top \mathbf{X}) \leq 4\sigma\sqrt{d}.$$

Moreover, for any $\delta > 0$, with probability $1 - \delta$, it holds

$$\|\mathbf{X}\| \leq 4\sigma\sqrt{d} + 2\sigma\sqrt{2\log(1/\delta)}.$$

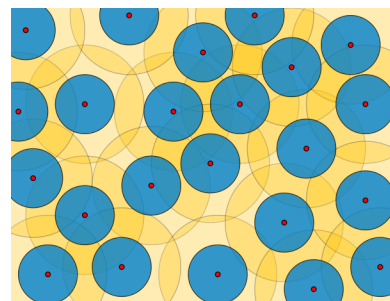


Figure 6.2: ϵ -cover of a square. Here in addition the points are at least ϵ apart from each other. Source: Wikipedia.

¹ The procedure needs to stop by compactness of \mathbb{B}_2 . If not, we would have a sequence (x_n) in \mathbb{B}_2 , which then would need to have a convergent subsequence. This would contradict the assumption that the elements of this sequence are all at least ϵ apart.

In this definition we do not require the components of \mathbf{X} to be independent

Proof. Let \mathcal{N} be a $1/2$ -cover of \mathbb{B}_2 that satisfies $|\mathcal{N}| \leq 6^d$. Next, observe that for every $\mathbf{u} \in \mathbb{B}_2$, there exists $\mathbf{y} \in \mathcal{N}$ and \mathbf{h} such that $\|\mathbf{h}\| \leq 1/2$ and $\mathbf{u} = \mathbf{y} + \mathbf{h}$. Therefore,

$$\max_{\mathbf{u} \in \mathbb{B}_2} \mathbf{u}^\top \mathbf{X} \leq \max_{\mathbf{y} \in \mathcal{N}} \mathbf{y}^\top \mathbf{X} + \max_{\mathbf{h} \in \frac{1}{2}\mathbb{B}_2} \mathbf{h}^\top \mathbf{X} = \max_{\mathbf{y} \in \mathcal{N}} \mathbf{y}^\top \mathbf{X} + \frac{1}{2} \max_{\mathbf{h} \in \mathbb{B}_2} \mathbf{h}^\top \mathbf{X}.$$

Therefore,

$$\max_{\mathbf{u} \in \mathbb{B}_2} \mathbf{u}^\top \mathbf{X} \leq 2 \max_{\mathbf{y} \in \mathcal{N}} \mathbf{y}^\top \mathbf{X}, \quad (6.3)$$

and so, using Proposition 6.1.1,

$$\mathbb{E} \max_{\mathbf{u} \in \mathbb{B}_2} \mathbf{u}^\top \mathbf{X} \leq 2 \mathbb{E} \max_{\mathbf{y} \in \mathcal{N}} \mathbf{y}^\top \mathbf{X} \leq 2\sigma \sqrt{2 \log |\mathcal{N}|} \leq 4\sigma \sqrt{d}.$$

The bound with high probability then follows because, using (6.3),

$$\mathbb{P}(\max_{\mathbf{u} \in \mathbb{B}_2} \mathbf{u}^\top \mathbf{X} \geq t) \leq \mathbb{P}(2 \max_{\mathbf{y} \in \mathcal{N}} \mathbf{y}^\top \mathbf{X} \geq t) \leq |\mathcal{N}| e^{-t^2/(8\sigma^2)} \leq 6^d e^{-t^2/(8\sigma^2)}.$$

Plugging $t = 4\sigma \sqrt{d} + 2\sigma \sqrt{2 \log(1/\delta)}$ we verify the bound. \square

← Exercise 6.5.1

6.2 Rademacher complexity and bounds on suprema

In this section we are again concerned with the variable

$$Z_n = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right| = \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}.$$

In Theorem 5.3.9, we showed that if \mathcal{F} is uniformly b -bounded, we get

$$\mathbb{P}(Z_n - \mathbb{E}Z_n \geq t) \leq \exp\left\{-\frac{nt^2}{2b^2}\right\}. \quad (6.4)$$

By setting $\delta = \exp\left\{-\frac{nt^2}{2b^2}\right\}$, we get $t = \sqrt{\frac{2b^2}{n} \log\left(\frac{1}{\delta}\right)}$. In other words, with probability at least $1 - \delta$, we have

$$0 \leq Z_n \leq \mathbb{E}Z_n + \sqrt{\frac{2b^2}{n} \log\left(\frac{1}{\delta}\right)}. \quad (6.5)$$

The second term on the right diminishes to 0 at the order $1/\sqrt{n}$. Since $\mathbb{E}Z$ is unknown, we need to bound it. The three main tools to construct such bounds are: symmetrization, discretization, and chaining. We now discuss them in more detail.

Definition 6.2.1 (Rademacher Complexity). *For a fixed collection $x_1^n := \{x_1, x_2, \dots, x_n\}$, $x_i \in \mathcal{X}$, consider*

$$\mathcal{F}(x_1^n) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}\} \subseteq \mathbb{R}^n. \quad (6.6)$$

The empirical Rademacher complexity is defined as

$$\mathcal{R}_n(\mathcal{F}(x_1^n)/n) := \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right].$$

Taking expectation w.r.t. X_1, X_2, \dots, X_n , we obtain the Rademacher complexity as

$$\mathcal{R}_n(\mathcal{F}) := \mathbb{E}_X \mathcal{R}_n(\mathcal{F}(X_1^n)/n).$$

Consider the following basic bound.

Theorem 6.2.2 (Symmetrization). *For any class of measurable functions \mathcal{F} , we have*

$$\mathbb{E} \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F}).$$

Proof. Suppose X'_1, \dots, X'_n are independent copies of X_1, \dots, X_n . As a result, we have

$$\begin{aligned} \mathbb{E} \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} &= \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X_i)) \right| \\ &= \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E}f(X'_i)) \right| \\ &\quad (\text{consider the convex function } g(y_1, \dots, y_n) = \sup_f |\frac{1}{n} \sum_i (f(x_i) - y_i)|) \\ &\stackrel{\text{Jensen}}{\leq} \mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right|. \end{aligned}$$

We have $f(X_i) - f(X'_i) \stackrel{d}{=} \varepsilon_i (f(X_i) - f(X'_i))$ for all i . This is why we call this the symmetrization argument. Thus we obtain

$$\begin{aligned} \mathbb{E} \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} &= \mathbb{E} \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(X_i) - f(X'_i)) \right| \right) \\ &\stackrel{\triangle}{\leq} 2\mathbb{E} \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \right) \\ &= 2\mathcal{R}_n(\mathcal{F}). \end{aligned}$$

□

As a corollary, using (6.5), we get the following important result.

Proposition 6.2.3. *For any uniformly b -bounded class of functions \mathcal{F} , any positive integer $n \geq 1$ and any scalar $t \geq 0$, we have*

$$\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F}) + t$$

with probability at least $1 - \exp(-\frac{nt^2}{2b^2})$. Consequently, as long as $\mathcal{R}_n(\mathcal{F}) = o(1)$, we have $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \rightarrow 0$ almost surely (i.e. \mathcal{F} is a Glivenko-Cantelli class for any \mathbb{P}).

Equivalently, $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 2\mathcal{R}_n(\mathcal{F}) + b\sqrt{\frac{2}{n} \log(\frac{1}{\delta})}$ with probability $\geq 1 - \delta$.

Proof. Let $Z_n := \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$. The first part follows directly from the earlier result and Theorem 5.3.9. In particular, convergence in probability of $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$ to zero follows easily². To get stronger almost sure convergence we employ the Borel-Cantelli lemma. By basic calculus, we have³

$$\{Z_n \rightarrow 0\} = \bigcap_{N=1}^{\infty} \bigcup_{n=1}^{\infty} \bigcap_{m \geq n} \{Z_m \leq \frac{1}{N}\}. \quad (6.7)$$

To prove almost sure convergence, we show that the complement of the set in (6.7) has measure zero. For a fixed N let $E_n = \{Z_n > \frac{1}{N}\}$. If n is large enough then $2\mathcal{R}_n(\mathcal{F}) < \frac{1}{2N}$ and so, for such n , taking $t = \frac{1}{2N}$ we get $\mathbb{P}(E_n) = \mathbb{P}(Z_n > \frac{1}{N}) \leq e^{-\frac{n}{2b^2N^2}}$. This implies that

$$\sum_{n \geq 1} \mathbb{P}(E_n) < \infty.$$

By the Borel-Cantelli lemma⁴, for every $N \geq 1$, $\mathbb{P}(\bigcup_{n=1}^{\infty} \bigcap_{m \geq n} \{Z_m \leq \frac{1}{N}\}) = 1$. Thus, in (6.7), we have a countable collection of measure 1 events. This implies that their intersection has measure 1 too, proving $\mathbb{P}(Z_n \rightarrow 0) = 1$. \square

To bound $\mathcal{R}_n(\mathcal{F})$, we first fix $x_1, \dots, x_n \in \mathcal{X}$ and bound the empirical Rademacher complexity $\mathcal{R}_n(\mathcal{F}(x_1^n)/n)$. Second, if the upper bound for $\mathcal{R}_n(\mathcal{F}(x_1^n)/n)$ does not depend on x_1, \dots, x_n , then it automatically becomes an upper bound for $\mathcal{R}_n(\mathcal{F})$.

We now state a simple upper bound for $\mathcal{R}_n(A) = \mathbb{E}[\sup_{a \in A} \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i]$ if $A \in \mathbb{R}^n$ is a set with finite elements. As a corollary from Proposition 6.1.1 we get the following bound.

Proposition 6.2.4 (Discretization). *Suppose A is a finite subset of \mathbb{R}^n with cardinality of $|A|$. Then*

$$\mathcal{R}_n(A) = \mathbb{E} \max_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i \right| \leq \max_{a \in A} \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2} \sqrt{\frac{2 \log(2|A|)}{n}}.$$

In Section 6.3 we show that the finite case already leads to deep and interesting results. Later, in Section 6.4 we show more general techniques.

6.3 Polynomial discrimination and VC dimension

In Proposition 6.2.4 we showed how to bound the Rademacher complexity for a finite set A . In this section we exploit this result.

Definition 6.3.1 (Boolean Class). *We say \mathcal{F} is a Boolean (function) class if $\forall f \in \mathcal{F}$ and $\forall x \in \mathcal{X}$, $f(x) \in \{0, 1\}$.*

² Recall that $Z_n \xrightarrow{p} 0$ if $\forall \epsilon > 0 \mathbb{P}(Z_n > \epsilon) \rightarrow 0$. Moreover, $Z_n \xrightarrow{a.s.} 0$ if $\mathbb{P}(Z_n \rightarrow 0) = 1$.

³ Recall $x_n \rightarrow 0$ if and only if $\forall N \geq 1 \exists n \geq 1$ s.t. $\forall m \geq n |x_m| \leq \frac{1}{N}$. Also, $\bigcup_{\alpha \in A} E_{\alpha} = \{x : \exists \alpha \text{ s.t. } x \in E_{\alpha}\}$ and $\bigcap_{\alpha \in A} E_{\alpha} = \{x : \forall \alpha \text{ s.t. } x \in E_{\alpha}\}$, which is why quantifier statements easily translate to set operations as in (6.7).

⁴ B-C states that if $\sum_{n \geq 1} \mathbb{P}(E_n) < \infty$ then $\mathbb{P}(\bigcap_{n \geq 1} \bigcup_{m \geq n} E_m) = 1$.

← Exercise 6.5.2

Recall that if ε_i are independent 1-sub-Gaussian then $\frac{1}{n} a^{\top} \varepsilon$ is $\frac{1}{n} \|a\|$ -sub-Gaussian.

An important example is given by the binary classification discussed in Section 4.1.1 and the uniform empirical process.

Example 6.3.2 (Binary Classification). *Consider a pair of random objects (X, Y) having some joint distribution where $X \in \mathcal{X}, Y \in \{0, +1\}$. A classifier is a function $g : \mathcal{X} \mapsto \{0, +1\}$. The error of the classifier is given by*

$$L(g) := \mathbb{P}(g(X) \neq Y) = \mathbb{E}\mathbb{1}(g(X) \neq Y).$$

Thus, in binary classification problems, the functions of interest are of the form $\mathbb{1}(g(X) \neq Y)$, and $\{\mathbb{1}(g(X) \neq Y) : g \in G\}$ is a Boolean class for any set of classifiers G .

Example 6.3.3 (Uniform Empirical Process). *When deriving the asymptotic law of the empirical CFD, the function class of interest $\mathcal{F} = \{\mathbb{1}_{(-\infty, t]}(\cdot) : t \in \mathbb{R}\}$ is a Boolean class.*

Now, fix a Boolean class \mathcal{F} and points $\{x_1, \dots, x_n\}$ in \mathcal{X} . Then $\mathcal{F}(x_1^n)$ defined in (6.6) is a finite set, contained in $\{0, 1\}^n$, whose cardinality is then at most 2^n . Applying Proposition 6.2.4, to $A = \mathcal{F}(x_1^n)/n$ we obtain

$$\mathcal{R}_n(\mathcal{F}(x_1^n)/n) = \mathbb{E} \max_{f \in \mathcal{F}} \frac{1}{n} \left| \sum_i \varepsilon_i f(x_i) \right| \leq \sqrt{\frac{2 \log(2|\mathcal{F}(x_1^n)|)}{n}}, \quad (6.8)$$

since the first term in Proposition 6.2.4 equals to $\sup_{f \in \mathcal{F}} \sqrt{\sum_{i=1}^n f^2(x_i)}/n$, which is less than or equal to 1.

Of course, if $|\mathcal{F}(x_1^n)| \approx 2^n$ then the bound in (6.8) is not very interesting. It becomes interesting, for example, when the cardinality of the function class grows only as a polynomial function of n , as formalized below.

Definition 6.3.4 (Polynomial Discrimination). *The Boolean class \mathcal{F} is said to have polynomial discrimination if there exists a polynomial $\rho(\cdot)$ such that for every positive integer n and every set of n points $\{x_1, \dots, x_n\}$, the set $\mathcal{F}(x_1^n)$ satisfies*

$$|\mathcal{F}(x_1^n)| \leq \rho(n).$$

The significance of this property is that, together with inequality (6.8), it provides a straightforward approach to controlling the Rademacher complexity, $\mathcal{R}_n(\mathcal{F}) = o(1)$. By Proposition 6.2.3 and (6.8), we get the following result.

Proposition 6.3.5. *If a Boolean class \mathcal{F} has polynomial discrimination then it is Glivenko-Cantelli class for any \mathbb{P} .⁵*

⁵ c.f. Definition 4.1.2

The first important example is given by the empirical process.

Proposition 6.3.6. *The Boolean class $\mathcal{F} = \{\mathbb{1}\{x \leq t\} : t \in \mathbb{R}\}$ satisfies $|\mathcal{F}(x_1^n)| \leq n + 1$. In particular, it has polynomial discrimination and so \mathcal{F} is a Glivenko-Cantelli class for any \mathbb{P} .*

But how does one check if a given Boolean class \mathcal{F} has polynomial discrimination in general? One of the most popular approaches is to use the Vapnik Chervonenkis dimension, or the VC dimension for short. We have the following lemma.

Definition 6.3.7 (VC Dimension). *The VC dimension $\text{VC}(\mathcal{F})$ of a class of Boolean functions \mathcal{F} on \mathcal{X} is defined as the **largest integer** D such that there exists a finite subset $\{x_1, \dots, x_D\}$ of \mathcal{X} satisfying $\mathcal{F}(x_1^D) = \{0, 1\}^D$.*

Note that it always holds that $\mathcal{F}(x_1^D) \subseteq \{0, 1\}^D$

Definition 6.3.8 (Shattering). *A finite subset $\{x_1, \dots, x_m\} \subset \mathcal{X}$ is said to be shattered by Boolean class \mathcal{F} if $\mathcal{F}(x_1^m) = \{0, 1\}^m$, i.e., $|\mathcal{F}(x_1^m)| = 2^m$.*

Remark 6.3.9. *By the definitions of VC dimension and shattering, we know the VC dimension is the largest integer n for which there is a n -point set $\{x_1, \dots, x_n\}$ which can be shattered by \mathcal{F} .*

Example 6.3.10. *Let \mathcal{F} be the set of indicator functions of the form $\mathbb{1}_{(-\infty, t]}(x)$ on \mathbb{R} . Then $\text{VC}(\mathcal{F}) = 1$. Indeed, if $x_1 < x_2$ then there is no $f \in \mathcal{F}$ for which $(f(x_1), f(x_2)) = (1, 0)$. By the same argument, no bigger set of points in \mathbb{R} can be shattered.*

The next lemma shows that any class with finite VC dimension has polynomial discrimination with degree of at most the VC dimension.

Lemma 6.3.11 (Sauer's Lemma). *Suppose $\text{VC}(\mathcal{F}) = D$, then for every $n \geq 1$ and every collection $\{x_1, \dots, x_n\}$ of n points, we have*

$$|\mathcal{F}(x_1^n)| \leq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{D}.$$

If $k > n$, we have $\binom{n}{k} = 0$, and if $D \leq n$, we have

$$|\mathcal{F}(x_1^n)| \leq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{D} \leq \left(\frac{en}{D}\right)^D.$$

Proof. See Proposition 4.18 of ⁶ for the detailed proof. □

⁶ Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019

As an immediate corollary, we get the following result.

Proposition 6.3.12. *Suppose \mathcal{F} is a Boolean class and $\text{VC}(\mathcal{F}) \leq D$. Then for $n \geq D$, we have*

$$\mathcal{R}_n(\mathcal{F}) \leq 2\sqrt{\frac{D}{n} \log\left(\frac{en}{D}\right)}$$

and so

$$\mathbb{E}\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 4\sqrt{\frac{D}{n} \log\left(\frac{en}{D}\right)}.$$

Proof. Combined with VC lemma and our previous upper bound (6.8), we can obtain that for $n \geq D$

$$\mathcal{R}_n(\mathcal{F}) \leq \sqrt{\frac{2 \log 2 + 2 \log \left(\frac{en}{D}\right)^D}{n}} \leq \sqrt{\frac{2(D+1) \log \left(\frac{en}{D}\right)}{n}} \leq 2\sqrt{\frac{D \log \left(\frac{en}{D}\right)}{n}},$$

where the second inequality follows because, for $n \geq D$, we have $2 \leq \left(\frac{en}{D}\right)$ and the last uses that $D \geq 1$. \square

6.3.1 Binary classification revisited

Consider again the classification problem discussed in Section 4.1.1. For a fixed family of classifiers \mathcal{C} , we study $\sup_{g \in \mathcal{C}} |R_n(g) - R(g)|$, where $R(g) = \mathbb{P}(g(X) \neq Y)$ and $R_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(X_i) \neq Y_i)$. The study of this quantity is well motivated by (4.2) and (4.3). This corresponds to $\|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}}$, where \mathcal{F} is taken to be the class of all functions $\mathbb{1}(g(x) \neq y)$ as g varies over \mathcal{C} (this class is uniformly 1-bounded), where the data are (X_i, Y_i) instead of X_i . Using the bounded differences inequality, Proposition 6.2.3, and Proposition 6.3.12, we get that with probability $\geq 1 - \delta$

$$\sup_{g \in \mathcal{C}} |R_n(g) - R(g)| = 2\sqrt{\frac{\text{VC}(\mathcal{F})}{n} \log\left(\frac{en}{\text{VC}(\mathcal{F})}\right)} + \sqrt{\frac{2}{n} \log\left(\frac{1}{\delta}\right)} \leq 2\sqrt{\text{VC}(\mathcal{F}) \frac{\log(n) + 1}{n}} + \sqrt{\frac{2}{n} \log\left(\frac{1}{\delta}\right)}.$$

To make these bounds efficient in a wide variety of examples, we formulate a bunch of results.

Lemma 6.3.13. *In the setting above, $\text{VC}(\mathcal{F}) \leq \text{VC}(\mathcal{C})$.*

Proof. We want to show that if \mathcal{F} can shatter $(x_1, y_1), \dots, (x_n, y_n)$, then \mathcal{C} can shatter x_1, \dots, x_n . For this, let $\eta_1, \dots, \eta_n \in \{0, 1\}$. We need to obtain $g \in \mathcal{C}$ such that $g(x_i) = \eta_i$ for $i = 1, \dots, n$. Define

$$\delta_i := \eta_i \mathbb{1}\{y_i = 0\} + (1 - \eta_i) \mathbb{1}\{y_i = 1\}.$$

As \mathcal{F} can shatter $(x_1, y_1), \dots, (x_n, y_n)$, there exists $f \in \mathcal{F}$, $f(x, y) = \mathbb{1}(g(x) \neq y)$ for some $g \in \mathcal{C}$ such that $f(x_i, y_i) = \delta_i$, $i = 1, \dots, n$. Then $g(x_i) = \eta_i$, for $i = 1, \dots, n$. This proves that \mathcal{C} shatters x_1, \dots, x_n . \square

When the logistic function is used for this classification task, the classifier is given by

$$g(x) = \mathbb{1}\{h(x) \geq \frac{1}{2}\}, \quad h(x) = \frac{1}{1 + e^{-x^\top \beta}}$$

or equivalently

$$g(x) = \mathbb{1}\{x^\top \beta \geq 0\}, \quad \beta \in \mathbb{R}^d.$$

Proposition 6.3.14 (Finite-Dimensional Vector Spaces). *Let \mathcal{V} be a d -dimensional vector space of real-valued functions on \mathcal{X} .⁷ Let $\mathcal{C} := \{\mathbb{1}(f \geq 0) : f \in \mathcal{V}\}$. Then $\text{VC}(\mathcal{C}) \leq d$.*

⁷ For example, $f = x^\top \beta$ with β running through \mathbb{R}^d .

Proof. For any fixed collection $\{x_1, \dots, x_{d+1}\}$ of $d+1$ points. Consider $\mathcal{T} = \{(f(x_1), \dots, f(x_{d+1})) : f \in \mathcal{V}\}$, then \mathcal{T} is a linear subspace of \mathbb{R}^{d+1} with dimension at most d . Therefore, there exists a nonzero vector $y \in \mathbb{R}^{d+1}$ such that

$$\sum_{i=1}^{d+1} y_i f(x_i) = 0, \quad \text{for all } f \in \mathcal{V}. \quad (6.9)$$

Without loss of generality, we assume that at least one y_k is positive, $k \in \{1, 2, \dots, d+1\}$ (If not, we can let $y' = -y$, then y' satisfies both (6.9) and the assumption). Now suppose there exists $\{x_1, \dots, x_{d+1}\}$ shattered by \mathcal{C} , meaning that each possible sign pattern of $(f(x_1), \dots, f(x_{d+1}))$ is possible. Then we can find a $f \in \mathcal{V}$ such that

$$\begin{aligned} f(x_i) &\geq 0, & \text{if } y_i &\leq 0, \\ f(x_i) &< 0, & \text{if } y_i &> 0. \end{aligned}$$

Thus we get $\sum_{i=1}^{d+1} y_i f(x_i) < 0$, which contradicts with (6.9). \square

Suppose now that a logistic classifier is considered and d is the dimension of the input space. Using these results, we conclude that with probability $\geq 1 - \delta$

$$\sup_{g \in \mathcal{C}} |R_n(g) - R(g)| \leq 2\sqrt{\frac{d(\log(n) + 1)}{n}} + \sqrt{\frac{2}{n} \log\left(\frac{1}{\delta}\right)}.$$

6.4 Chaining

The situation like described in Section 6.3 is very special. Typically we do not have that the functions $f \in \mathcal{F}$ all take values in a finite set. In general, to exploit the discretization argument in Proposition 6.2.4 we need to work harder.

Suppose now that $A \subseteq \mathbb{R}^n$ is arbitrary. Let $N = N(\delta)$ be the δ -covering number of A with respect to the norm $\|a\|_n := \|a\|/\sqrt{n}$. A naive approach for bounding $\mathcal{R}_n(A)$ is as follows. Fix a δ -covering \mathcal{N} of A and for every a denote by $\pi(a)$ the closest point in \mathcal{N} to a we get

$$\mathcal{R}_n(A) = \mathbb{E} \sup_{u \in A} \frac{1}{n} |\varepsilon^\top u| \leq \mathbb{E} [\max_{z \in \mathcal{N}} \frac{1}{n} |\varepsilon^\top z|] + \mathbb{E} [\sup_{u \in A} \frac{1}{n} |\varepsilon^\top (u - \pi(u))|] \leq \sqrt{\frac{2 \log(2N(\delta))}{n}} \sup_{u \in A} \frac{\|u\|}{\sqrt{n}} + \delta,$$

where the last inequality follows by Proposition 6.2.4 and the Cauchy-Schwarz inequality. This inequality holds for every δ but $N(\delta)$ monotonically increases as δ decreases.

Get back to our problem of bounding the empirical Rademacher complexity $\mathcal{R}(\mathcal{F}(X_1^n)/n)$. In this section, we only assume that the class \mathcal{F} has a finite and integrable envelope function F : a function such that $|f(x)| \leq F(x) < \infty$, for every x and $f \in \mathcal{F}$.

Consider now the $L_2(\mathbb{P})$ -norm, $\|f\|_{\mathbb{P},2} = \sqrt{\mathbb{P}|f|^2}$. Given a subset \mathcal{F} of the $L_2(\mathbb{P})$ -space, we denote the δ -covering number by $N(\delta, \mathcal{F}, L_2(\mathbb{P}))$. Note that

$$\|f\|_{\mathbb{P}_n,2} = \sqrt{\frac{1}{n} \sum_{i=1}^n f^2(X_i)} \quad \text{and} \quad \sup_{f \in \mathcal{F}} \|f\|_{\mathbb{P}_n,2} \leq \|F\|_{\mathbb{P}_n,2},$$

which is the same as $\|a\|_n$ with $a = (f(X_1), \dots, f(X_n))$. In other words, $N(\delta)$ for $A = \mathcal{F}(X_1^n)/n$ above is the same as $N(\delta, \mathcal{F}, L_2(\mathbb{P}_n))$. Thus,

$$\mathcal{R}_n(\mathcal{F}(X_1^n)/n) \leq \sqrt{\frac{2 \log N(\delta, \mathcal{F}, L_2(\mathbb{P}_n))}{n}} \|F\|_{\mathbb{P}_n,2} + \delta,$$

which implies

$$\begin{aligned} \mathcal{R}_n(\mathcal{F}(X_1^n)/n) &\leq \left(\sqrt{\frac{2 \log N(\delta \|F\|_{\mathbb{P}_n,2}, \mathcal{F}, L_2(\mathbb{P}_n))}{n}} + \delta \right) \|F\|_{\mathbb{P}_n,2} \\ &\leq \left(\sqrt{\frac{2 \log \sup_Q N(\delta \|F\|_{Q,2}, \mathcal{F}, L_2(Q))}{n}} + \delta \right) \|F\|_{\mathbb{P}_n,2}, \end{aligned}$$

where the supremum is taken over all measures supported on a finite set with n points⁸. From this, we can conclude the following result, which we state without proof.

Theorem 6.4.1. *Let \mathcal{F} be a suitably measurable class of measurable functions with $\sup_Q N(\delta \|F\|_{L_2(Q)}, \mathcal{F}, L_2(Q)) < \infty$ for every $\delta > 0$. If $\mathbb{E}F(X) < \infty$ then \mathcal{F} is \mathbb{P} -Glivenko-Cantelli.*

A tighter bound on the Rademacher complexity can be obtained using chaining.

Theorem 6.4.2 (Chaining). *Suppose $0 \in A \subseteq \mathbb{R}^n$ and let $D := \sup_{a \in A} \|a\| / \sqrt{n}$. Then*

$$\mathcal{R}_n(A) \leq \frac{16}{\sqrt{n}} \int_0^{D/2} \sqrt{\log N(\delta)} d\delta.$$

Proof. Let A_m be a $D/2^m$ -covering of A with $|A_m| = N(D/2^m)$. We have $A_0 = \{0\}$. For every $m \geq 0$ let $\pi_m(a)$ denote the point in A_m that is the closest to a . Note that

$$\frac{1}{n} \varepsilon^\top a = \sum_{m=0}^{\infty} \frac{1}{n} \varepsilon^\top (\pi_{m+1}(a) - \pi_m(a)).$$

⁸ Note that this gives a trivial bound on the Rademacher complexity $\mathbb{E}\mathcal{R}(\mathcal{F}(X_1^n)/n)$ because, by Jensen's inequality $\mathbb{E}\|F\|_{\mathbb{P}_n,2} \leq \|F\|_{\mathbb{P},2}$ and this is the only term in the bound above that depends on X_1^n .

By Proposition 6.2.4, we have

$$\begin{aligned}
\mathbb{E}[\max_{a \in A} \frac{1}{n} \varepsilon^\top (\pi_{m+1}(a) - \pi_m(a))] &\leq \sqrt{\frac{2 \log(2|A_m| \cdot |A_{m+1}|)}{n}} \max_{a \in A} \frac{\|\pi_{m+1}(a) - \pi_m(a)\|}{\sqrt{n}} \\
&= \sqrt{\frac{2 \log(2|A_{m+1}|^2)}{n}} \max_{a \in A} \frac{\|\pi_{m+1}(a) - a + a - \pi_m(a)\|}{\sqrt{n}} \\
&\leq \frac{4}{2^m} D \sqrt{\frac{\log(|A_{m+1}|)}{n}} \\
&\leq \frac{4}{2^m} D \sqrt{\frac{\log(N(D/2^{m+1}))}{n}}.
\end{aligned}$$

This assures that

$$\begin{aligned}
\mathbb{E}[\sup_{a \in A} \frac{1}{n} |\varepsilon^\top a|] &\leq \sum_{m=0}^{\infty} \mathbb{E}[\max_{a \in A} \frac{1}{n} |\varepsilon^\top (\pi_{m+1}(a) - \pi_m(a))|] \\
&\leq 4 \sum_{m=0}^{\infty} \frac{D}{2^m} \sqrt{\frac{\log N(D/2^{m+1})}{n}} \\
&\leq 16 \sum_{m=0}^{\infty} \frac{D}{2^{m+2}} \sqrt{\frac{\log N(D/2^{m+1})}{n}} \\
&\leq 16 \int_0^{D/2} \sqrt{\frac{\log N(\delta)}{n}} d\delta
\end{aligned}$$

as claimed. \square

Theorem 6.4.3. Suppose there exists a function F such that $|f(x)| \leq F(x)$ for all $f \in \mathcal{F}$. Then

$$\mathbb{E} \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq 8 \sqrt{\frac{\mathbb{E}(F(X)^2)}{n}} \int_0^1 \sqrt{\log \sup_Q N(\delta \|F\|_{L^2(Q)}, \mathcal{F}, L^2(Q))} d\delta,$$

where the supremum is taken over all measures supported on a finite set with n points.

Proof. Consider a fixed sample $(X_i)_{i=1}^n$ with the underlying distribution \mathbb{P}_n . Note that

$$\sup_{a \in \mathcal{F}(X_1^n)} \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2} = \sup_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum_{i=1}^n f(X_i)^2} \leq \sqrt{\frac{1}{n} \sum_{i=1}^n F(X_i)^2} = \sqrt{\mathbb{P}_n F^2}$$

By Theorem 6.4.2, conditionally on $(X_i)_{i=1}^n$

$$\begin{aligned}
\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(X_i) \right| \middle| (X_i)_{i=1}^n \right] &\leq \frac{16}{\sqrt{n}} \int_0^{\sqrt{\mathbb{P}_n F^2}/2} \sqrt{\log N(\delta; \mathcal{F}, L^2(\mathbb{P}_n))} d\delta \\
&= \frac{8}{\sqrt{n}} \sqrt{\mathbb{P}_n F^2} \int_0^1 \sqrt{\log N(\frac{\delta}{2} \sqrt{\mathbb{P}_n F^2}; \mathcal{F}, L^2(\mathbb{P}_n))} d\delta \\
&\leq \frac{8}{\sqrt{n}} \sqrt{\mathbb{P}_n F^2} \int_0^1 \sup_Q \sqrt{\log N(\frac{\delta}{2} \|F\|_{L^2(Q)}; \mathcal{F}, L^2(Q))} d\delta.
\end{aligned}$$

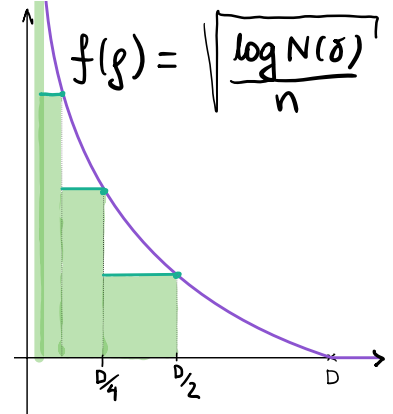


Figure 6.3: Illustration of the last bound in the proof of Theorem 6.4.2.

To get the desired bound we now take the expectation with respect to $(X_i)_{i=1}^n$. By the Jensen inequality we get

$$\mathbb{E} \sqrt{\frac{\mathbb{P}_n F^2}{n}} \leq \sqrt{\frac{\mathbb{E} F^2(X)}{n}},$$

which gives the desired result. \square

We conclude stating the following result that shows that Proposition 6.3.12 can be strengthened.

Proposition 6.4.4. *Suppose \mathcal{F} is a Boolean class and $\text{VC}(\mathcal{F}) \leq D$. Then for any $\delta > 0$, we have*

$$\sup_Q N(\delta, \mathcal{F}, L^2(Q)) \leq \left(\frac{c_1}{\delta}\right)^{c_2 D},$$

where c_1, c_2 are some positive (universal) constants, and the supremum is over all probability distributions over \mathcal{X} . Consequently,

$$\mathbb{E} \|\mathbb{P}_n - \mathbb{P}\|_{\mathcal{F}} \leq c \sqrt{\frac{D}{n}}$$

for some universal constant c .

6.5 Exercises

Exercise 6.5.1. Let $X \in \mathbb{R}^{n \times d}$ be a random matrix with i.i.d. entries that are σ -sub-Gaussian. Denote by S^{n-1} the unit sphere in \mathbb{R}^n and by \mathbb{B}_2^n the corresponding unit ball.

- (i) Show that $u^T X v$ is a σ -sub-Gaussian random variable for any $u \in \mathbb{B}_2^n$, $v \in \mathbb{B}_2^d$.
- (ii) The operator norm $\|X\|$ is defined as $\|X\| := \sup_{v \neq 0} \frac{\|Xv\|}{\|v\|}$. Show that

$$\|X\| = \max_{u \in S^{n-1}} \max_{v \in S^{d-1}} u^T X v = \max_{u \in \mathbb{B}_2^n} \max_{v \in \mathbb{B}_2^d} u^T X v.$$

- (iii) Using (i) and (ii), show that there exists a constant $C > 0$ such that $\mathbb{E} \|X\| \leq C(\sqrt{n} + \sqrt{d})$.

Exercise 6.5.2. Prove Proposition 6.2.4.

7

Applications in statistics (1-2 weeks)

Please read Chapter 1 in ¹ to build up some intuition for why high-dimensional statistics is important in general. In here we only focus on a basic illustration for the results in the previous chapter. For more details, see the lecture notes of Philippe Rigollet², on which I based most of this chapter.

¹ Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019

² <https://math.mit.edu/~rigollet/PDFs/RigNotes17.pdf>

7.1 Sub-Gaussian sequence model with sparsity

Consider the Gaussian sequence model, that is, let

$$X_i = \mu_i + \epsilon_i, \quad i = 1, \dots, d, \quad (7.1)$$

where $\epsilon_i \sim N(0, \sigma^2)$. This is the same model as we considered in Section 2.5 in our discussion of the Stein's paradox. Occasionally, we will also relax the assumption of Gaussianity of ϵ in which case we refer to it as a sub-Gaussian sequence model.

Given a random sample $X^{(1)}, \dots, X^{(n)} \in \mathbb{R}^d$ from model (7.1), we can estimate μ by the sample mean $Y = \bar{X}_n = \frac{1}{n} \sum_{j=1}^n X^{(j)}$. Note that $\mathbb{E}Y_i = \mu_i$ and $\text{var}(Y_i) = \sigma^2/n$ and thus, equivalently, we can consider a Gaussian sequence model

$$Y_i = \mu_i + \epsilon_i, \quad i = 1, \dots, d,$$

where $\text{var}(\epsilon_i) = \sigma^2/n$ or, more generally, a sub-Gaussian sequence model for which ϵ_i is $\frac{\sigma}{\sqrt{n}}$ -sub-Gaussian. We use the estimator $\hat{\mu} = (Y_1, \dots, Y_d)$ with the risk³

$$R(\mu, \hat{\mu}) = \mathbb{E} \left[\sum_{i=1}^d \epsilon_i^2 \right] \leq \frac{\sigma^2 d}{n}.$$

To think about this as a high dimensional problem, we let d, n be both large in which case $\hat{\mu}$ may not be a good estimator unless $n \gg d$.

In order to estimate μ in the high-dimensional setting we will require some additional structure on μ . In Section 2.5 we saw a simple

³ This formula already appeared in (2.11).

The inequality becomes equality in the Gaussian case.

example of how explicit bounds on the entries of μ can be exploited by an estimator of the form CY with C diagonal.

Here we assume that μ is sparse or most of its entries are small. In this case, a natural estimator is given by hard thresholding. Hard thresholding gives the estimator

$$\hat{\mu}_i^\tau = Y_i \mathbb{1}(|Y_i| \geq \tau) \quad \text{for all } i = 1, \dots, d.$$

Equivalently, the thresholding estimator is a solution to the problem:

$$\hat{\mu}^\tau = \arg \min_{a \in \mathbb{R}} \left\{ \|Y - a\|^2 + \tau^2 \sum_{i=1}^d \mathbb{1}\{a_i \neq 0\} \right\}.$$

We will try to obtain risk bounds for hard thresholding. Note that, by (6.1), with probability at least $1 - \delta$,

$$\max_i |\epsilon_i| \leq \sigma \sqrt{\frac{2 \log(2d/\delta)}{n}} =: \tau. \quad (7.2)$$

The consequences of this inequality are two-fold. First, if $\mu_i = 0$ then $|Y_i| = |\epsilon_i| \leq \tau$ with high probability. Thus, if we observe $|Y_i| \gg \tau$ then it must correspond to $\mu_i \neq 0$. Second, if $|Y_i| \leq \tau$ then μ_i cannot be too large because, by the triangle inequality,

$$|\mu_i| \leq |Y_i| + |\epsilon_i| \leq 2\tau.$$

Therefore, we loose at most 2τ by taking $\hat{\mu}_i^\tau = 0$.

Proposition 7.1.1. *Let $\hat{\mu}^{2\tau}$ be the hard thresholding estimator with threshold 2τ , where τ is defined in (7.2) Then, (7.2) implies that*

$$\|\hat{\mu}^{2\tau} - \mu\|^2 \leq 9 \sum_{i=1}^d \min\{\mu_i^2, \tau^2\}. \quad (7.3)$$

In particular, (7.3) holds with probability $\geq 1 - \delta^4$

Proof. We condition on the high probability event (7.2), which we call \mathcal{E} in this proof. Fix index i and note that:

1. If $|\mu_i| \leq \tau$ then (given \mathcal{E}) $|Y_i| \leq |\mu_i| + |Y_i - \mu_i| \leq 2\tau$ and so $\hat{\mu}_i^{2\tau} = 0$. In this case, $(\mu_i - \hat{\mu}_i^{2\tau})^2 = \mu_i^2$.
2. If $|\mu_i| > 3\tau$ then (given \mathcal{E}) $|Y_i| \geq |\mu_i| - |Y_i - \mu_i| > 2\tau$ and so $\hat{\mu}_i^{2\tau} = Y_i$. In this case $(\mu_i - \hat{\mu}_i^{2\tau})^2 = \epsilon_i^2 \leq \tau^2$.
3. If $\tau < |\mu_i| \leq 3\tau$, then

$$(\hat{\mu}_i^{2\tau} - \mu_i)^2 = (Y_i \mathbb{1}(|Y_i| \geq 2\tau) - \mu_i)^2 = \mu_i^2 \mathbb{1}(|Y_i| < 2\tau) + \epsilon_i^2 \mathbb{1}(|Y_i| \geq 2\tau) \leq \max\{\epsilon_i^2, \mu_i^2\} \leq 9\tau^2.$$

Putting these together, we see that (7.2) implies (7.3). \square

In each a_i the function is Y_i^2 if $a_i = 0$ and $(Y_i - a_i)^2 + \tau^2$ otherwise. The latter function is minimized at $a_i = Y_i$ with the optimal value τ^2 . Thus the optimum of the whole function depends on whether $|Y_i| \geq \tau$ or not.

Recall: ϵ_i is $\frac{\sigma}{\sqrt{n}}$ -sub-Gaussian.

⁴ Recall that (7.3) depends on δ through (7.2).

The next theorem shows how this can be exploited. It is convenient to have the following standard notation:

$$a_n \lesssim b_n \quad \text{means} \quad a_n \leq Cb_n \quad \text{for some } C > 0 \text{ and all } n.$$

Theorem 7.1.2. *Consider the same set-up as in Proposition 7.1.1. If (7.2) holds, then the following statements hold:*

(i) *If $\|\mu\|_0 = s$ then*

$$R(\mu, \hat{\mu}^{2\tau}) \leq 9s\tau^2 = 18\sigma^2 \frac{s \log(2d/\delta)}{n} \lesssim \sigma^2 \frac{s \log(d)}{n}. \quad (7.4)$$

Note that we do not bound the risk but the conditional risk. But here the event we condition on holds with high probability $(1 - \delta)$ so this is still informative.

(ii) *If $\min_{i \in \text{supp}(\mu)} |\mu_i| > 3\tau$, then*

$$\text{supp}(\hat{\mu}^{2\tau}) = \text{supp}(\mu).$$

Thus we get some guarantees on support recovery.

Equation (7.4) shows that in the sparse setting $\hat{\mu}^\tau$ may be consistent as long as $\frac{s \log(d)}{n} \rightarrow 0$, where s is the number of non-zero entries of the vector μ . Quite surprisingly, this may happen even if d is exponentially larger than n .

A natural question is whether sparsity is required for such a high-dimensional consistency result. There are several examples showing that this is not the case. To get a flavour of these results suppose that $\|\mu\|_1 \leq R$ for some $R \in \mathbb{R}$. Let $\tau = \sigma \sqrt{\frac{\log(2d/\delta)}{n}}$ as defined in (7.2). If $\|\mu\|_1 \leq R$, then the number of i such that $|\mu_i| > \tau$ is at most R/τ . By Theorem 7.1.1, under the event \mathcal{E} defined in (7.2),

$$R(\hat{\mu}^{2\tau}, \mu) \leq 9 \sum_{i=1}^d \min\{\mu_i^2, \tau^2\}$$

Thus, we can use the previous (conditional) risk bound to obtain:

$$\begin{aligned} R(\hat{\mu}^{2\tau}, \mu) &\leq 9 \sum_{i=1}^d \min\{\mu_i^2, \tau^2\} \\ &= 9 \sum_{i: |\mu_i| > \tau} \tau^2 + 9 \sum_{i: |\mu_i| \leq \tau} \mu_i^2 \\ &\leq \frac{9R}{\tau} \tau^2 + 9 \sum_{i: |\mu_i| \leq \tau} \mu_i^2 \\ &\leq 9R\tau + 9\tau \sum_{i: |\mu_i| \leq \tau} |\mu_i| \\ &\leq 18R\tau = 18R\sigma \sqrt{\frac{\log(2d/\delta)}{n}}. \end{aligned}$$

Notice that the rate of convergence is different from the sparse case, roughly behaving as $1/\sqrt{n}$ instead of $1/n$. Nevertheless, this is still a strong, nice and surprising result.

7.2 Fixed design linear regression

Consider the standard linear regression problem

$$\mathbf{y} = \mathbf{X}\theta^* + \epsilon,$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is a fixed matrix, $\theta^* \in \mathbb{R}^d$ is a fixed vector and ϵ is a random vector such that each ϵ_i is an independent random variable with $\mathbb{E}\epsilon_i = 0$, $\text{var}(\epsilon_i) = \sigma^2$. Given the data \mathbf{y} , the least squares estimator of θ^* is obtained by minimizing

$$M_n(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \theta)^2 = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\theta\|^2 \quad (7.5)$$

with respect to $\theta \in \mathbb{R}^d$. If $\text{rank}(\mathbf{X}) = d$ (in part. $n \geq d$) then the optimum exists and is unique. It is given by the well known formula

$$\bar{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

If $n < d$ then the optimum is attained over an affine subspace of positive dimension and so it is not unique. It is then customary to fix

$$\bar{\theta} = (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{y},$$

where A^+ is the pseudo-inverse of A . In this case $\bar{\theta}$ minimizes the norm over all the optimal points.

The quality of an estimator $\hat{\theta}$ is normally analysed in this context based on the mean square error $\mathbb{E}\|\hat{\theta} - \theta^*\|^2$. For simplicity we will first study the mean squared prediction error

$$\text{MSE}(\mathbf{X}\hat{\theta}) = \frac{1}{n} \|\mathbf{X}(\hat{\theta} - \theta^*)\|^2 = (\hat{\theta} - \theta^*)^\top \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X}\right) (\hat{\theta} - \theta^*). \quad (7.6)$$

As we will see later, properly choosing the design also ensures that $\|\hat{\theta} - \theta^*\|^2$ is small as long as $\text{MSE}(\mathbf{X}\hat{\theta})$ is small.

We are now going to prove our first result on the finite sample performance of the least squares estimator for fixed design.

Theorem 7.2.1. *Assume that the linear model holds with ϵ which is σ -sub-Gaussian. Then the least squares estimator $\bar{\theta}$ satisfies*

$$\mathbb{E}(\text{MSE}(\mathbf{X}\bar{\theta})) = \frac{1}{n} \mathbb{E}\|\mathbf{X}(\bar{\theta} - \theta^*)\|^2 \leq 4\sigma^2 \frac{r}{n} \lesssim \sigma^2 \frac{r}{n},$$

where $r = \text{rank}(\mathbf{X}^\top \mathbf{X})$. Moreover, for any $\delta \geq 0$, with probability at least $1 - \delta$, it holds

$$\text{MSE}(\mathbf{X}\bar{\theta}) \leq 32\sigma^2 \frac{2r + \log(1/\delta)}{n}$$

Proof. By definition

$$\|\mathbf{y} - \mathbf{X}\bar{\theta}\|^2 \leq \|\mathbf{y} - \mathbf{X}\theta^*\|^2 = \|\epsilon\|^2. \quad (7.7)$$

More generally, we could consider a feature map $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ and the corresponding feature matrix $\Psi \in \mathbb{R}^{n \times d}$ whose rows are obtained by applying ψ to the rows of \mathbf{X} . The corresponding model would be

$$\mathbf{y} = \Psi\theta^* + \epsilon.$$

Moreover, denoting $\bar{\Delta} = \bar{\theta} - \theta^*$,

$$\|\mathbf{y} - \mathbf{X}\bar{\theta}\|^2 = \|\mathbf{X}\theta^* + \epsilon - \mathbf{X}\bar{\theta}\|^2 = \|\mathbf{X}\bar{\Delta}\|^2 - 2\epsilon^\top \mathbf{X}\bar{\Delta} + \|\epsilon\|^2.$$

This, together with (7.7), allows us to conclude that

$$\|\mathbf{X}\bar{\Delta}\|^2 \leq 2\epsilon^\top \mathbf{X}\bar{\Delta}. \quad (7.8)$$

Using the Cauchy-Schwarz inequality directly for the right-hand side in (7.8) would be wasteful in the case when $r \ll d$. Instead, let $V = [\mathbf{v}_1 \cdots \mathbf{v}_r] \in \mathbb{R}^{n \times r}$ be a matrix whose r columns form an orthonormal basis of the column span of \mathbf{X} . In particular, there exists \mathbf{u} such that $\mathbf{X}\bar{\Delta} = V\mathbf{u}$ and, denoting $\tilde{\epsilon} = V^\top \epsilon$, we have

$$\epsilon^\top \mathbf{X}\bar{\Delta} = \epsilon^\top V\mathbf{u} = \tilde{\epsilon}^\top \mathbf{u} \leq \|\tilde{\epsilon}\| \|\mathbf{u}\| = \|\tilde{\epsilon}\| \|\mathbf{X}\bar{\Delta}\|,$$

where we used the fact that $\|V\mathbf{u}\| = \|\mathbf{u}\|$ as $V^\top V = I_r$. Using this inequality in (7.8) gives that

$$\|\mathbf{X}\bar{\Delta}\| \leq 2\|\tilde{\epsilon}\| \quad (7.9)$$

and so

$$\mathbb{E}\|\mathbf{X}(\bar{\theta} - \theta^*)\|^2 \leq 4\mathbb{E}\|\tilde{\epsilon}\|^2 = 4 \sum_{i=1}^r \mathbb{E}\tilde{\epsilon}_i^2 \leq 4r\sigma^2,$$

where the last inequality follows from the fact that each entry $\tilde{\epsilon}_i = \mathbf{v}_i^\top \epsilon$ is zero-mean σ -sub-Gaussian and Exercise 5.5.6. This concludes the proof of the bound on $\mathbb{E}(\text{MSE}(\mathbf{X}\bar{\theta}))$.

For the second statement note that, by (7.9),

$$\mathbb{P}(\text{MSE}(\mathbf{X}\bar{\theta}) \geq t) = \mathbb{P}(\|\mathbf{X}(\bar{\theta} - \theta^*)\|^2 \geq nt) \leq \mathbb{P}(\|\tilde{\epsilon}\|^2 \geq nt/4) = \mathbb{P}(\|\tilde{\epsilon}\| \geq \sqrt{nt}/2).$$

Since $\|\tilde{\epsilon}\| = \sup_{\mathbf{u} \in \mathbb{B}^2} \mathbf{u}^\top \tilde{\epsilon}$ we can use the last inequality in the proof of Theorem 6.1.6 to conclude that

$$\mathbb{P}(\|\tilde{\epsilon}\| \geq \sqrt{nt}/2) \leq 6^r e^{-\frac{nt}{32\sigma^2}}.$$

Taking $t = 32\sigma^2 \frac{2r + \log(1/\delta)}{n}$ we verify that

$$\mathbb{P}(\text{MSE}(\mathbf{X}\bar{\theta}) \geq t) \leq \delta.$$

□

If $r = d \leq n$ then bounds on prediction errors give bounds on $\|\bar{\theta} - \theta^*\|^2$. In this case $B = \frac{1}{n}\mathbf{X}^\top \mathbf{X}$ has rank d and, using (7.6), we get

$$\|\bar{\theta} - \theta^*\|^2 \leq \frac{\text{MSE}(\mathbf{X}\bar{\theta})}{\gamma_{\min}(B)}, \quad (7.10)$$

where $\gamma_{\min}(B)$ is the minimal eigenvalue of B . Theorem 7.2.1 can be therefore used to bound $\|\bar{\theta} - \theta^*\|^2$ directly. By contrast, in the high dimensional setting B is not positive definite and we need more structure.

7.3 Constrained least squares estimator

Let $K \subseteq \mathbb{R}^d$ be a symmetric convex set⁵. If we knew a priori that $\theta^* \in K$, we may prefer a constrained least squares estimator $\bar{\theta}_K$ defined by

$$\bar{\theta}_K \in \arg \min_{\theta \in K} \|\mathbf{y} - \mathbf{X}\theta\|^2.$$

The equivalent of inequality (7.8) still holds, that is, $\|\mathbf{X}(\bar{\theta}_K - \theta^*)\|^2 \leq 2\epsilon^\top \mathbf{X}(\bar{\theta}_K - \theta^*)$. Further,

$$\|\mathbf{X}(\bar{\theta}_K - \theta^*)\|^2 \leq 2\epsilon^\top \mathbf{X}(\bar{\theta}_K - \theta^*) \leq 2 \sup_{\theta \in K-K} \epsilon^\top \mathbf{X}\theta,$$

where $K - K = \{x - y : x, y \in K\}$. It is easy to see that since K is symmetric and convex $K - K = 2K$ so that

$$2 \sup_{\theta \in K-K} \epsilon^\top \mathbf{X}\theta = 4 \sup_{v \in \mathbf{X}K} \epsilon^\top v$$

where $\mathbf{X}K = \{\mathbf{X}\theta : \theta \in K\} \subseteq \mathbb{R}^n$. This is the measure of the size of $\mathbf{X}K$. If $\epsilon \sim N(0, I_d)$, the expected value of the above supremum is called the **Gaussian width** of $\mathbf{X}K$.

ℓ_1 constrained least squares Assume here that $K = \mathbb{B}_1$ is the unit ℓ_1 ball of \mathbb{R}^d . Recall that it has exactly $2d$ vertices $\pm e_1, \dots, \pm e_d$, where e_i is the i -th canonical unit vector. It implies that the set $\mathbf{X}\mathbb{B}_1$ is also a polytope with at most $2d$ vertices that are contained in the set $\{-\mathbf{X}_1, \mathbf{X}_1, \dots, \mathbf{X}_d, \mathbf{X}_d\}$, where \mathbf{X}_i is the i -th column of \mathbf{X} .

Theorem 7.3.1. *Suppose $\theta^* \in \mathbb{B}_1$. Moreover, assume the conditions of Theorem 7.2.1 and that the columns of \mathbf{X} are normalized so that $\max_i \|\mathbf{X}_i\| \leq \sqrt{n}$. Then the constrained least squares estimator $\bar{\theta}_{\mathbb{B}_1}$ satisfies*

$$\mathbb{E}[\text{MSE}(\mathbf{X}\bar{\theta}_{\mathbb{B}_1})] = \frac{1}{n} \mathbb{E} \|\mathbf{X}(\bar{\theta}_{\mathbb{B}_1} - \theta^*)\|^2 \lesssim \sigma \sqrt{\frac{\log d}{n}}.$$

Moreover, for any $\delta \in (0, 1)$, with probability $1 - \delta$, it holds

$$\text{MSE}(\mathbf{X}\bar{\theta}_{\mathbb{B}_1}) \lesssim \sigma \sqrt{\frac{\log(d/\delta)}{n}}.$$

Proof. From the considerations preceding the theorem, we got that

$$\|\mathbf{X}(\bar{\theta}_{\mathbb{B}_1} - \theta^*)\|^2 \leq 4 \sup_{v \in \mathbf{X}\mathbb{B}_1} \epsilon^\top v.$$

Moreover, because $\mathbf{X}\mathbb{B}_1$ is a polytope, we have

$$\sup_{v \in \mathbf{X}\mathbb{B}_1} \epsilon^\top v = \max_{i=1, \dots, d} |\epsilon^\top \mathbf{X}_i|.$$

Since ϵ is σ -sub-Gaussian⁶, then for any column \mathbf{X}_i such that $\|\mathbf{X}_i\| \leq$

⁵ Symmetric means that $K = -K$.

← Exercise 7.6.1

⁶ Recall this multivariate definition given before Theorem 6.1.6.

\sqrt{n} , the random variable $\epsilon^\top \mathbf{X}_i$ is $(\sqrt{n}\sigma)$ -sub-Gaussian. Therefore, applying Proposition 6.1.1, we get

$$\mathbb{E} \sup_{v \in \mathbf{X}\mathbb{B}_1} \epsilon^\top v \leq \sigma \sqrt{2n \log(2d)},$$

which gives the claimed bound on $\mathbb{E}[\text{MSE}(\mathbf{X}\bar{\theta}_{\mathbb{B}_1})]$. Again by Proposition 6.1.1 we get that, for any $t > 0$

$$\mathbb{P}(\|\mathbf{X}(\bar{\theta}_{\mathbb{B}_1} - \theta^*)\|^2 \geq 4t) \leq \mathbb{P}(\sup_{v \in \mathbf{X}\mathbb{B}_1} \epsilon^\top v \geq t) \leq 2de^{-t^2/(2\sigma^2n)}.$$

From this we conclude⁷

$$\mathbb{P}(\text{MSE}(\mathbf{X}\hat{\theta}) \geq t) \leq 2de^{-nt^2/(32\sigma^2)}.$$

To conclude the proof, we find t such that

$$2de^{-nt^2/(32\sigma^2)} \leq \delta \Leftrightarrow t^2 \geq 32\sigma^2 \frac{\log(2d)}{n} + 32\sigma^2 \frac{\log(1/\delta)}{n}.$$

□

Note that the proof of Theorem 7.2.1 also applies to $\bar{\theta}_{\mathbb{B}_1}$ so that $\bar{\theta}_{\mathbb{B}_1}$ benefits from the best of both rates

$$\mathbb{E}[\text{MSE}(\mathbf{X}\bar{\theta}_{\mathbb{B}_1})] \lesssim \min\left\{\frac{r}{n}, \sqrt{\frac{\log d}{n}}\right\}.$$

This is called an elbow effect. The elbow takes place around $r \simeq \sqrt{n \log d}$.

⁷ Recall from (7.6) that $\text{MSE}(\mathbf{X}\hat{\theta}) = \frac{1}{n} \|\mathbf{X}(\hat{\theta} - \theta^*)\|^2$.

← Exercise 7.6.2

7.4 LASSO regression

The focus of this section is on the situation when the true parameter vector θ^* is sparse. In this case, even if $n < d$ it may be possible to control the error of recovering θ^* . Building upon (7.5), we consider the regularized estimator

$$\hat{\theta} := \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\theta\|^2 + \lambda \|\theta\|_1 \right\}, \tag{7.11}$$

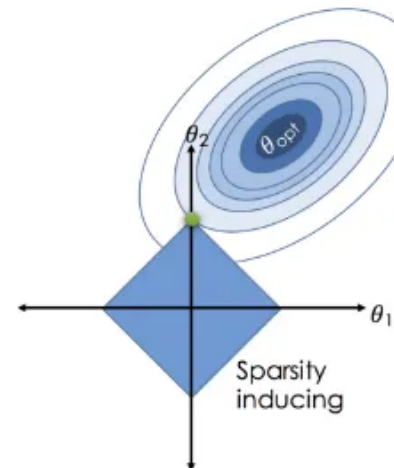
where $\lambda \geq 0$ is a fixed penalty parameter. This is called the LASSO regression problem.

The LASSO regression problem is a convex problem but it is not differentiable. By standard convex duality theory assures that, for any given $\lambda \geq 0$, the LASSO problem (7.11) is equivalent to the constrained problem

$$\text{minimize } \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\theta\|^2 \quad \text{subject to } \|\theta\|_1 \leq R \tag{7.12}$$

for some $R \geq 0$. This observation aids some intuition behind what the LASSO estimator is actually doing (contemplate the side figure).

Although the optimum is not given in a closed form, there is a simple numerical algorithm that can be used for optimization. The algorithm relies on the observation that if $d = 1$ then the optimum is given in a closed form. Using this fact, we can run a coordinate descent algorithm where at each step we update one coordinate of θ keeping the other coordinates fixed.



There are three main problems related with the analysis of the LASSO estimator. One focuses on $\|\hat{\theta} - \theta^*\|$ to establish high-dimensional consistency. The other focuses on the prediction performance $\|\mathbf{X}(\hat{\theta} - \theta^*)\|$ as analysed for the least squares estimator in the previous section. Finally, we can study how the LASSO regression in recovering the true support of θ^* . Below we briefly focus on the first and the last.

7.4.1 High-dimensional consistency

Recall that $M_n(\theta) = \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\theta\|^2$. In the discussion surrounding (7.10) we argued that if the spectrum of $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ (and so the Hessian of $M_n(\theta)$) is bounded away from zero, a bound on the predictive error provide bounds on $\text{MSE}(\hat{\theta})$. Considering the dual problem (7.12) with R satisfying $\|\theta^*\|_1 \leq R$, we could conclude something similar for the LASSO estimator. In the high-dimensional setting $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ is not positive definite. However, using the fact that θ^* is sparse, we need enough of curvature of M_n only in some directions.

The Hessian matrix $\nabla^2 M_n(\theta) = \frac{1}{n} \mathbf{X}^T \mathbf{X}$ is positive definite with minimal eigenvalue at least κ if

$$\frac{1}{n} \Delta^T \mathbf{X}^T \mathbf{X} \Delta = \frac{1}{n} \|\mathbf{X}\Delta\|^2 \geq \kappa \|\Delta\|^2 \quad \text{for all } \Delta \in \mathbb{R}^d.$$

As we said, if $n < d$, there is no $\kappa > 0$ for which this condition holds. Instead, for any $S \subseteq \{1, \dots, d\}$ and $\alpha > 0$ define

$$\mathbf{C}_\alpha(S) := \{\Delta \in \mathbb{R}^d : \|\Delta\|_1 \leq (1 + \alpha) \|\Delta_S\|_1\}.$$

We say that the matrix \mathbf{X} satisfies the **restricted eigenvalue (RE) condition** over S with parameters (κ, α) if

$$\frac{1}{n} \|\mathbf{X}\Delta\|^2 \geq \kappa \|\Delta\|^2 \quad \text{for all } \Delta \in \mathbf{C}_\alpha(S).$$

Suppose that:

- (A1) The vector θ^* is supported on a subset $S \subseteq \{1, \dots, d\}$ with $|S| = s$.
- (A2) The design matrix satisfies the restricted eigenvalue condition over S with parameters $(\kappa, 3)$.

Theorem 7.4.1. *Under assumptions (A1) and (A2) any solution of (7.11) for $\lambda_n \geq 2 \|\frac{\mathbf{X}^T \epsilon}{n}\|_\infty$ we have*

$$\|\hat{\theta}_n - \theta^*\| \leq \frac{3}{\kappa} \sqrt{s} \lambda_n.$$

Large portion of $\|\Delta\|_1$ is due to Δ_S . In other words $\|\Delta_{S^c}\|_1 \leq \alpha \|\Delta_S\|_1$

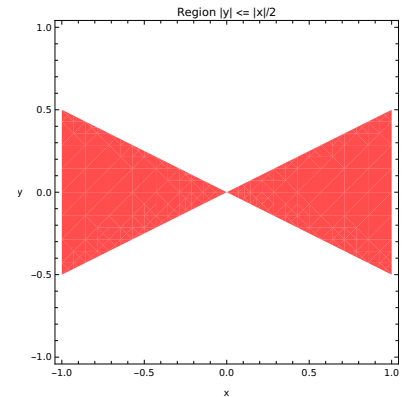


Figure 7.1: Suppose that $d = 2$, $\theta = (x, y)$, $S = \{1\}$, and $\alpha = \frac{1}{4}$. Then $\mathbf{C}_\alpha(S) = \{(x, y) : |y| \leq \frac{1}{2}|x|\}$.

Proof. We first show that, if $\lambda_n \geq 2\|\frac{\mathbf{X}^\top \epsilon}{n}\|_\infty$ then the error $\widehat{\Delta} = \widehat{\theta} - \theta^*$ belongs to $\mathbf{C}_3(S)$. Let $L(\theta; \lambda_n) = \frac{1}{2n}\|\mathbf{y} - \mathbf{X}\theta\|^2 + \lambda_n\|\theta\|_1$. We have

$$L(\widehat{\theta}; \lambda_n) \leq L(\theta^*; \lambda_n) = \frac{1}{2n}\|\epsilon\|^2 + \lambda_n\|\theta^*\|_1.$$

With $\widehat{\Delta} = \widehat{\theta} - \theta^*$ we get

$$\|\mathbf{y} - \mathbf{X}\widehat{\theta}\|^2 = \|\mathbf{y} - \mathbf{X}\theta^* - \mathbf{X}\widehat{\Delta}\|^2 = \|\epsilon - \mathbf{X}\widehat{\Delta}\|^2.$$

We can use it in the previous expression to conclude

$$0 \leq \frac{1}{2n}\|\mathbf{X}\widehat{\Delta}\|^2 \leq \frac{1}{n}\epsilon^\top \mathbf{X}\widehat{\Delta} + \lambda_n(\|\theta^*\|_1 - \|\widehat{\theta}\|_1). \quad (7.13)$$

Since θ^* is S -sparse, we can write

$$\|\theta^*\|_1 - \|\widehat{\theta}\|_1 = \|\theta_S^*\|_1 - \|\theta_S^* + \widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1.$$

Substituting this into (7.13) gives

$$\begin{aligned} 0 &\leq \frac{1}{n}\|\mathbf{X}\widehat{\Delta}\|^2 \leq \frac{2}{n}\epsilon^\top \mathbf{X}\widehat{\Delta} + 2\lambda_n(\|\theta_S^*\|_1 - \|\theta_S^* + \widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1) \\ &\leq 2\left\|\frac{\mathbf{X}^\top \epsilon}{n}\right\|_\infty \|\widehat{\Delta}\|_1 + 2\lambda_n(\|\widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1) \\ &\leq \lambda_n(3\|\widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1), \end{aligned}$$

where the last inequality follows from the choice of λ_n . The fact that $3\|\widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1 \geq 0$ establishes that $\widehat{\Delta} \in \mathbf{C}_3(S)$ so that the RE condition can be applied. Doing so, we conclude that

$$\kappa\|\widehat{\Delta}\|^2 \leq \frac{1}{n}\|\mathbf{X}\widehat{\Delta}\|^2 \leq \lambda_n(3\|\widehat{\Delta}_S\|_1 - \|\widehat{\Delta}_{S^c}\|_1) \leq 3\lambda_n\|\widehat{\Delta}_S\|_1.$$

Since $\|\widehat{\Delta}_S\|_1 \leq \sqrt{s}\|\widehat{\Delta}_S\| \leq \sqrt{s}\|\widehat{\Delta}\|$, the conclusion follows. \square

We now show how this result can be applied in the classical linear Gaussian model for which the noise vector ϵ has i.i.d. $N(0, \sigma^2)$ entries. More generally, the same calculation applies when ϵ is σ -sub-Gaussian.

In addition, we assume that \mathbf{X} satisfies the RE condition and that it is C -column normalized, meaning that⁸

$$\max_j \|\mathbf{X}_j\| \leq C\sqrt{n},$$

where \mathbf{X}_j denotes the columns of \mathbf{X} . With this set-up, the random variable $\|\frac{1}{n}\mathbf{X}^\top \epsilon\|_\infty$ corresponds to the absolute maximum of d zero mean Gaussian variables, each with variance at most $C^2\sigma^2/n$. Consequently from standard sub-Gaussian tail bounds in Proposition 6.1.1

$$\mathbb{P}\left(\left\|\frac{1}{n}\mathbf{X}^\top \epsilon\right\|_\infty \geq t\right) \leq 2de^{-\frac{nt^2}{2C^2\sigma^2}}.$$

⁸ Exercise 7.6.3 partially motivates this assumption.

← Exercise 7.6.3

Plugging, $t = C\sigma \left(\sqrt{\frac{2\log d}{n}} + \delta \right)$ we easily verify that

$$\mathbb{P} \left(\left\| \frac{1}{n} \mathbf{X}^T \epsilon \right\|_\infty \geq C\sigma \left(\sqrt{\frac{2\log d}{n}} + \delta \right) \right) \leq 2e^{-n\delta^2/2} \quad \text{for all } \delta > 0.$$

Consequently, if we set

$$\lambda_n = 2C\sigma \left(\sqrt{\frac{2\log d}{n}} + \delta \right)$$

then Theorem 7.4.1 implies that

$$\|\hat{\theta} - \theta^*\| \leq \frac{6C\sigma}{\kappa} \sqrt{s} \left\{ \sqrt{\frac{2\log d}{n}} + \delta \right\}$$

with probability at least $1 - 2e^{-n\delta^2/2}$.

7.4.2 Model recovery consistency

Here we again focus on the deterministic design. For variable selection consistency the restricted eigenvalue condition is replaced by a closely related condition.

(A3) **Lower eigenvalue:** The smallest eigenvalue of the sample covariance submatrix indexed by S is bounded below

$$\gamma_{\min} \left(\frac{1}{n} \mathbf{X}_S^T \mathbf{X}_S \right) \geq c_{\min} > 0.$$

(A4) **Mutual incoherence:** There exists $\alpha \in [0, 1)$ such that

$$\max_{j \in S^c} \|(\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T X_j\|_1 \leq \alpha.$$

Note that condition (A3) is rather mild and it is required to get identifiability even if S is known in advance. The second condition is more subtle and roughly it says that no variables in S^c are too correlated with the support variables. This is form of orthogonality ($(\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T X_j$ is a projection of X_j on the span of \mathbf{X}_S), which is actually unlikely to hold in big datasets.

Recall that for a vector $w \in \mathbb{R}^d$, $\|w\|_\infty = \max_i |w_i|$. If $W \in \mathbb{R}^{k \times l}$ is a matrix then $\|W\|_\infty = \max_{i=1, \dots, k} \|W_i\|_1$, where W_i is the i -th row of W .

Theorem 7.4.2. Consider an S -sparse linear regression model for which the design matrix satisfies conditions (A3) and (A4). If

$$\lambda_n \geq \frac{2}{1-\alpha} \left\| \frac{1}{n} \mathbf{X}_{S^c}^T \left(I_n - \mathbf{X}_S (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T \right) \epsilon \right\|_\infty,$$

then $\hat{\theta}$ has the following properties:

- (a) *Uniqueness: There is a unique optimal solution $\widehat{\theta}$.*
- (b) *No false inclusion: This solution has its support set \widehat{S} contained within the true support set S .*
- (c) *ℓ_∞ -bounds: The error $\widehat{\Delta}$ satisfies*

$$\|\widehat{\Delta}_S\|_\infty \leq \|(\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \epsilon\|_\infty + \|(\frac{1}{n} \mathbf{X}_S^\top \mathbf{X}_S)^{-1}\|_\infty \lambda_n =: B(\lambda_n, \mathbf{X}).$$

- (d) *No false exclusion: The lasso includes all $i \in S$ such that $|\theta_i^*| > B(\lambda_n, \mathbf{X})$, and hence is variable selection consistent if $\min_{i \in S} |\theta_i^*| > B(\lambda_n, \mathbf{X})$.*

Corollary 7.4.3. *Consider the S -sparse linear model based on a noise vector ϵ with zero-mean i.i.d. σ -sub-Gaussian entries, and a deterministic design matrix \mathbf{X} that satisfies (A3) and (A4), as well as the C -column normalization condition. The LASSO estimator with*

$$\lambda_n = \frac{2C\sigma}{1-\alpha} \left\{ \sqrt{\frac{2 \log(d-s)}{n}} + \delta \right\}$$

for some $\delta > 0$. Then $\widehat{\theta}$ is unique with its support contained within S (no type I errors), and satisfies the ℓ_∞ -error bound

$$\|\widehat{\Delta}_S\|_\infty \leq \frac{\sigma}{\sqrt{c_{\min}}} \left\{ \sqrt{\frac{2 \log s}{n}} + \delta \right\} + \|(\frac{1}{n} \mathbf{X}_S^\top \mathbf{X}_S)^{-1}\|_\infty \lambda_n$$

all with probability at least $1 - 4e^{-n\delta^2/2}$.

Both results are left without a proof, see Chapter 7 in ⁹ for details.

⁹ Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, Cambridge, 2019

7.5 Random matrices

7.5.1 Spectral norm of sub-Gaussian random matrices

Let $A \in \mathbb{R}^{m \times n}$, then $\|A\| = \max_{\|\mathbf{x}\| \leq 1} \|A\mathbf{x}\|$ is called the operator norm of A . Using the variational characterization of the norm, we get

$$\|A\| = \max_{\|\mathbf{x}\| \leq 1, \|\mathbf{y}\| \leq 1} \mathbf{y}^\top A \mathbf{x}.$$

Lemma 7.5.1. *Let \mathcal{N} be an ϵ -covering of \mathbb{B}_2^n . Then*

$$\max_{\mathbf{x} \in \mathcal{N}} \|A\mathbf{x}\| \leq \|A\| \leq \frac{1}{1-\epsilon} \max_{\mathbf{x} \in \mathcal{N}} \|A\mathbf{x}\|.$$

If, in addition, \mathcal{M} is an ϵ -covering of \mathbb{B}_2^m then

$$\max_{\mathbf{x} \in \mathcal{N}, \mathbf{y} \in \mathcal{M}} \mathbf{y}^\top A \mathbf{x} \leq \|A\| \leq \frac{1}{1-2\epsilon} \max_{\mathbf{x} \in \mathcal{N}, \mathbf{y} \in \mathcal{M}} \mathbf{y}^\top A \mathbf{x}.$$

Proof. For the first statement – we used this argument already in the proof of Theorem 6.1.6. For the second statement, use the fact that $\|Ax\| = \max_{\|y\| \leq 1} y^\top Ax$ and that $(1 - \epsilon)^2 \geq 1 - 2\epsilon$. \square

Proposition 7.5.2. *Let $A \in \mathbb{R}^{m \times n}$ be a random matrix with independent, mean zero, σ -sub-Gaussian entries. Then, for all $t \geq 0$*

$$\|A\| \leq 5\sigma(\sqrt{m} + \sqrt{n} + t)$$

with probability $\geq 1 - e^{-t^2}$.

Proof. First, consider $\frac{1}{4}$ -coverings \mathcal{M}, \mathcal{N} of \mathbb{B}_2^m and \mathbb{B}_2^n . By Lemma 6.1.4, we can assume $|\mathcal{M}| \leq 12^m$ and $|\mathcal{N}| \leq 12^n$. By Lemma 7.5.1, $\|A\| \leq 2 \max_{y \in \mathcal{M}, x \in \mathcal{N}} y^\top Ax$. By Lemma 5.2.4, each $y^\top Ax$ is σ -sub-Gaussian. By the union bound we conclude

$$\mathbb{P}(\|A\| \geq u) \leq \mathbb{P}\left(\max_{y \in \mathcal{M}, x \in \mathcal{N}} y^\top Ax \geq \frac{u}{2}\right) \leq 12^{n+m} e^{-\frac{u^2}{8\sigma^2}}.$$

Take $u^* = 5\sigma(\sqrt{m} + \sqrt{n} + t)$. Then

$$\mathbb{P}(\|A\| \geq u^*) \leq 12^{m+n} e^{-\frac{25}{8}(m+n+t^2)} \leq e^{-t^2}.$$

\square

We note also that $\mathbb{E}\|A\| \leq 5\sigma(\sqrt{m} + \sqrt{n})$. This follows from Proposition 6.1.1. Indeed,

$$\mathbb{E}\|A\| \leq 2\mathbb{E}\left[\max_{y \in \mathcal{M}, x \in \mathcal{N}} y^\top Ax\right] \leq 2\sigma\sqrt{2 \log 12^{m+n}} \leq 5\sigma\sqrt{m+n},$$

from which the conclusion follows.

7.5.2 Recovering communities in the stochastic block model

We consider a basic version of the stochastic block model. This is a model for random graphs on n nodes, which are divided into two equal-sized communities. The model comes with two parameters $0 < q < p < 1$. Each edge $i - j$ appears independently with probability p if i, j are in the same community and q if they are in two different communities.

We represent a graph with an adjacency matrix $A \in \{0, 1\}^{n \times n}$, where

$$\mathbb{P}(A_{ij} = 1) = \mathbb{E}A_{ij} = \begin{cases} p & \text{if } i, j \text{ lie in the same community,} \\ q & \text{if } i, j \text{ lie in different communities.} \end{cases}$$

If the vertices are ordered so that the first $n/2$ belong to the first community then $\mathbb{E}A$ has a simple block structure

$$\mathbb{E}A = \begin{bmatrix} p\mathbf{1}\mathbf{1}^\top & q\mathbf{1}\mathbf{1}^\top \\ q\mathbf{1}\mathbf{1}^\top & p\mathbf{1}\mathbf{1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} p & q \\ q & p \end{bmatrix} \begin{bmatrix} \mathbf{1}^\top & \mathbf{0}^\top \\ \mathbf{0}^\top & \mathbf{1}^\top \end{bmatrix}$$

This section is adapted from:

Roman Vershynin. *High-dimensional probability: An Introduction with Applications in Data Science*, volume 47. Cambridge University Press, 2018

where $\mathbf{0}, \mathbf{1} \in \mathbb{R}^{n/2}$ are vector of zeros and ones. Denote

$$u_1 = \begin{bmatrix} \mathbf{1} \\ \mathbf{1} \end{bmatrix}, \quad u_2 = \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix}, \quad v_1 = \frac{1}{\sqrt{n}}u_1, \quad v_2 = \frac{1}{\sqrt{n}}u_2,$$

where v_1, v_2 are the normalized versions of u_1, u_2 with $\|v_1\| = \|v_2\| = 1$. It is then easy to check that

$$\mathbb{E}Av_1 = n \frac{p+q}{2} v_1 \quad \text{and} \quad \mathbb{E}Av_2 = n \frac{p-q}{2} v_2.$$

In other words, u_1, u_2 are the eigenvectors of $\mathbb{E}A$ with eigenvalues $\lambda_1 = n \frac{p+q}{2}$ and $\lambda_2 = n \frac{p-q}{2}$ respectively. Since $\mathbb{E}A$ has rank 2, all other eigenvalues are zero.

Consider now situation when we observe A but the community structure is unknown. If we observed $\mathbb{E}A$, we could compute the eigenvector corresponding to the second largest eigenvalue and then we could assign vertices to communities according to whether the corresponding entry in this eigenvector is positive or negative. A natural question is what happens if we do the same on the observed matrix A .

By the Weyl's theorem we know that for any two symmetric matrices $\max_i |\lambda_i(S) - \lambda_i(T)| \leq \|S - T\|$. We have a similar result for eigenvectors.

Theorem 7.5.3 (Davis-Kahan). *If S, T are symmetric $n \times n$ matrices. Fix i and suppose $\min_{j \neq i} |\lambda_i(S) - \lambda_j(S)| = \delta > 0$ then*

$$\sin(\angle\{v_i(S), v_i(T)\}) \leq \frac{2\|S - T\|}{\delta}.$$

Note that the sine of an angle being small means that the angle between $v_i(S)$ and $v_i(T)$ is either close to 0 or close to π . We can always find $\theta \in \{-1, 1\}$ such that the angle between $v_i(S)$ and $\theta v_i(T)$ is close to zero.

Proposition 7.5.4. *With the same assumption as in Theorem 7.5.3, there exists $\theta \in \{-1, 1\}$ such that*

$$\|v_i(S) - \theta v_i(T)\| \leq \frac{2^{3/2}\|S - T\|}{\delta}.$$

Proof. Let $u = v_i(S)$, $v = \theta v_i(T)$, where θ is fixed so that $u^\top v \geq 0$. Note that $\|u\| = \|v\| = 1$ and $\|u - v\|^2 = 2(1 - u^\top v)$. Moreover, using the formula $\cos(\angle(u, v)) = u^\top v$ we get

$$\sin^2(\angle(u, v)) = 1 - (u^\top v)^2 \geq 1 - u^\top v = \frac{1}{2}\|u - v\|^2.$$

By Theorem 7.5.3, we then conclude

$$\|u - v\|^2 \leq 2 \sin^2(\angle(u, v)) \leq \frac{8\|S - T\|}{\delta^2},$$

from which the result follows. \square

We are going to use this result with $S = \mathbb{E}A$ and $T = A$. To use the David-Kahan theorem, we first check that λ_2 (the second eigenvalue of $\mathbb{E}A$) is well-separated from the rest of the spectrum, that is, from λ_1 and 0. We have

$$\delta = \min\{\lambda_1 - \lambda_2, \lambda_2 - 0\} = n \min\left\{\frac{p-q}{2}, q\right\} =: n\mu.$$

In SBM A_{ij} 's are independent. Since $|A_{ij} - \mathbb{E}A_{ij}| \leq 1$, it follows that $A_{ij} - \mathbb{E}A_{ij}$ are independent 1-sub-Gaussian random variables. It follows that

$$\|A - \mathbb{E}A\| \leq 10(\sqrt{n} + t)$$

with probability $\geq 1 - e^{-t^2}$. By Proposition 7.5.4, there exists $\theta \in \{-1, 1\}$ such that

$$\|v_i(\mathbb{E}A) - \theta v_i(A)\| \leq \frac{2^{3/2}\|A - \mathbb{E}A\|}{\delta} \leq \frac{5 \cdot 2^{5/2}}{\mu} \left(\frac{t}{n} + \frac{1}{\sqrt{n}}\right)$$

with probability $\geq 1 - e^{-t^2}$. Take $t = \sqrt{n}$ to conclude that

$$\|v_i(\mathbb{E}A) - \theta v_i(A)\| \leq \frac{C}{\mu\sqrt{n}} \quad (C := 5 \cdot 2^{7/2})$$

with probability $\geq 1 - e^{-n}$.

Let $u_i = \sqrt{n}v_i$. Note that the entries of $u_i(\mathbb{E}A)$ are ± 1 , so if signs of $u_i(\mathbb{E}A)$ and $\theta u_i(A)$ do not agree, they contribute > 1 to $\|u_i(\mathbb{E}A) - \theta u_i(A)\|^2$. We know that, with high probability,

$$\|u_i(\mathbb{E}A) - \theta u_i(A)\|^2 \leq \frac{C^2}{\mu^2}$$

and so there cannot be more than $\frac{C^2}{\mu^2}$ (a constant) of entries that contribute more than 1!

We conclude that, with high probability, we can correctly classify all but finite number of vertices. We formulate this as a theorem.

Theorem 7.5.5 (Spectral clustering for SBMs). *Let $A \sim \text{SBM}(n, p, q)$ with $p > q$ and $\min\{p - q, q\} =: \mu > 0$. Then with probability $\geq 1 - e^{-n}$ the spectral clustering algorithm identifies communities of A correctly up to C^2/μ^2 misspecified vertices.*

7.5.3 Covariance matrix estimation

The basic problem can be formulated as follows. Suppose $\mathbf{X} \in \mathbb{R}^{n \times d}$ such that each row \mathbf{x}_i is i.i.d. $N_d(0, \Sigma)$. The sample covariance matrix is defined as

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X}^\top \mathbf{X}.$$

This is a random positive semi-definite matrix whose expectation satisfies

$$\mathbb{E}\widehat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \mathbb{E}\mathbf{x}_i\mathbf{x}_i^\top = \Sigma.$$

The question we ask is how $\widehat{\Sigma}_n$ concentrates around Σ . Ideally, we would also like to relax the Gaussianity assumption.

Concentration can be measured in various norms. If we use the operator norm, the general idea is to use the following variational representation: If $Q \in \mathbb{S}^d$ then

$$\|Q\| = \sup_{\|v\|=1} |v^\top Q v|.$$

Thus

$$\|\widehat{\Sigma}_n - \Sigma\| = \sup_{\|v\|=1} |v^\top \widehat{\Sigma}_n v - v^\top \Sigma v| = \sup_{\|v\|=1} \left| \frac{1}{n} \sum_{i=1}^n (v^\top \mathbf{x}_i)^2 - v^\top \Sigma v \right|.$$

Let \mathcal{F} be the class of functions $f(\mathbf{x}) = (v^\top \mathbf{x})^2$ for $\|v\| = 1$. Since $\mathbb{E}f(\mathbf{x}) = v^\top \Sigma v$, we can rewrite

$$\|\widehat{\Sigma}_n - \Sigma\| = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) - \mathbb{E}f(\mathbf{X}) \right|.$$

Note that the spectrum of $\widehat{\Sigma}_n$ is directly related to the spectrum of \mathbf{X} . We have

$$\gamma_{\max}(\widehat{\Sigma}_n) = \gamma_{\max}\left(\frac{1}{n}\mathbf{X}^\top \mathbf{X}\right) = \left(\sigma_{\max}\left(\frac{1}{\sqrt{n}}\mathbf{X}\right)\right)^2$$

and

$$\gamma_{\min}(\widehat{\Sigma}_n) = \gamma_{\min}\left(\frac{1}{n}\mathbf{X}^\top \mathbf{X}\right) = \left(\sigma_{\min}\left(\frac{1}{\sqrt{n}}\mathbf{X}\right)\right)^2.$$

Since singular values are continuous functions of its matrix argument, we then expect that

$$\sigma_{\max}\left(\frac{1}{\sqrt{n}}\mathbf{X}\right) = \sqrt{\gamma_{\min}(\widehat{\Sigma}_n)} \approx \sqrt{\gamma_{\min}(\Sigma)} = \gamma_{\min}(\sqrt{\Sigma}).$$

Indeed, we have the following result.

Theorem 7.5.6 (The Gaussian case). *Suppose $\mathbf{X} \in \mathbb{R}^{n \times d}$ such that each row \mathbf{x}_i is i.i.d. from $N_d(0, \Sigma)$. Then for every $\delta > 0$*

$$\mathbb{P}\left(\sigma_{\max}\left(\frac{1}{\sqrt{n}}\mathbf{X}\right) \geq (1 + \delta)\gamma_{\max}(\sqrt{\Sigma}) + \sqrt{\frac{\text{tr}(\Sigma)}{n}}\right) \leq e^{-n\delta^2/2}$$

and, if $n \geq d$, we also have

$$\mathbb{P}\left(\sigma_{\min}\left(\frac{1}{\sqrt{n}}\mathbf{X}\right) \leq (1 - \delta)\gamma_{\min}(\sqrt{\Sigma}) - \sqrt{\frac{\text{tr}(\Sigma)}{n}}\right) \leq e^{-n\delta^2/2}$$

Proof. The proof, as many results of that form has two step. First we show that $\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X})$ concentrates around its expectation. Second, we provide bounds on this expectation that allow us to conclude the given formulas.

Step I: We can write $\mathbf{X} = \mathbf{W}\sqrt{\Sigma}$, where \mathbf{W} has i.i.d. standard normal rows. The first part can be shown as in Example 5.4.4. By Weyl's theorem $\mathbf{W} \mapsto \sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{W}\sqrt{\Sigma})$ is Lipschitz with parameter $L = \frac{1}{\sqrt{n}}\|\sqrt{\Sigma}\|$. By Theorem 5.4.1

$$\mathbb{P}(\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X}) \geq \mathbb{E}\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X}) + t) \leq e^{-nt^2/(2\|\Sigma\|)}.$$

Taking $t = \sqrt{\|\Sigma\|}\delta$,

$$\mathbb{P}(\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X}) \geq \mathbb{E}\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X}) + \sqrt{\|\Sigma\|}\delta) \leq e^{-n\delta^2/2}.$$

Preparing for the second part of the proof, note that if we could show that $\mathbb{E}(\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X})) \leq \sqrt{\|\Sigma\|} + \sqrt{\frac{\text{tr}(\Sigma)}{n}}$, we would be done.

Step II: We want to show that $\mathbb{E}(\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X})) \leq \sqrt{\|\Sigma\|} + \sqrt{\frac{\text{tr}(\Sigma)}{n}}$. We use the following variational characterization

$$\sigma_{\max}(\frac{1}{\sqrt{n}}\mathbf{X}) = \max_{\|\mathbf{u}\|=1} \max_{\|\mathbf{v}\|=1} \mathbf{u}^\top (\frac{1}{\sqrt{n}}\mathbf{X}) \mathbf{v}.$$

Let $S^{n-1} \subset \mathbb{R}^n$ be the set of vectors with unit norm. Consider the zero-mean Gaussian process

$$Z_{\mathbf{u},\mathbf{v}} := \mathbf{u}^\top (\frac{1}{\sqrt{n}}\mathbf{X}) \mathbf{v} \quad \text{for } (\mathbf{u}, \mathbf{v}) \in S^{n-1} \times S^{d-1}.$$

Consider the induced metric on $S^{n-1} \times S^{d-1}$

$$d((\mathbf{u}, \mathbf{v}), (\mathbf{u}', \mathbf{v}')) := \sqrt{\mathbb{E}(Z_{\mathbf{u},\mathbf{v}} - Z_{\mathbf{u}',\mathbf{v}'})^2}.$$

We have

$$\mathbb{E}(Z_{\mathbf{u},\mathbf{v}} - Z_{\mathbf{u}',\mathbf{v}'})^2 \leq \|\Sigma\| \|\mathbf{u} - \mathbf{u}'\|^2 + \|\mathbf{v} - \mathbf{v}'\|_{\Sigma}^2.$$

We construct another Gaussian process $Y_{\mathbf{u},\mathbf{v}}$, whose covariance function is equal to this right-hand side. \square

7.6 Exercises

Exercise 7.6.1. Show that if $K \subset \mathbb{R}^d$ is convex and symmetric ($K = -K$) then $K - K = \{x - y : x, y \in K\}$ is equal to $2K = \{2x : x \in K\}$.

Exercise 7.6.2. Show that the proof of Theorem 7.2.1 applies to the estimator $\bar{\theta}_{\mathbb{B}_1}$.

Exercise 7.6.3. Suppose $\mathbf{X} \in \mathbb{R}^{n \times d}$ is a design matrix, whose rows are independent observations of a d -variate mean zero distribution for which all marginal distributions are σ -sub-Gaussian. Show that that, for every j , the j column \mathbf{X}_j of \mathbf{X} satisfies $\|\mathbf{X}_j\| \leq 4\sigma\sqrt{n} + t$ with probability $\geq 1 - e^{-t^2/(8\sigma^2)}$.

Part III

Some topics from the first semester

8

Classical large sample theory

This section will be mostly based on a few chapters from ¹. We will be relying on basic definitions and results on convergence in probability theory.

¹ A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998

8.1 Preliminaries

We consider vectors with values in $\mathcal{X} \subseteq \mathbb{R}^m$. The set \mathcal{X} forms a metric space with the induced metric $d(x, y) = \|x - y\|$ (but any other equivalent metric is fine). The inequality $x \leq y$ is meant coordinatewise. Many of the results discussed here can be generalized to arbitrary measure spaces with the underlying Borel measure (open sets are measurable).

Recall that a sequence of random vectors X_n converges in distribution to X , denoted $X_n \rightsquigarrow X$, if

$$\mathbb{P}(X_n \leq x) \rightarrow \mathbb{P}(X \leq x),$$

for every point x for which the limit CDF $x \mapsto \mathbb{P}(X \leq x)$ is continuous. Convergence in distribution is often called the weak convergence. Moreover, X_n converges to X in probability, denoted $X_n \xrightarrow{P} X$, if for all $\epsilon > 0$

$$\mathbb{P}(d(X_n, X) > \epsilon) \rightarrow 0.$$

You have studied various equivalent formulations of weak convergence. The one useful for us is given in item (iv) in the next lemma.

Lemma 8.1.1 (Portmanteau Lemma). *We have $X_n \rightsquigarrow X$ if and only if any of the following conditions holds:*

- (i) $\mathbb{E}f(X_n) \rightarrow \mathbb{E}f(X)$ for all bounded, continuous functions f ,
- (ii) $\mathbb{E}f(X_n) \rightarrow \mathbb{E}f(X)$ for all bounded, Lipschitz functions f ,
- (iii) $\liminf \mathbb{E}f(X_n) \geq \mathbb{E}f(X)$ for all nonnegative, continuous functions f ,

(iv) $\liminf \mathbb{P}(X_n \in U) \geq \mathbb{P}(X \in U)$ for every open set $U \subseteq \mathcal{X}$,

(v) $\limsup \mathbb{P}(X_n \in B) \leq \mathbb{P}(X \in B)$ for every closed set $B \subseteq \mathcal{X}$,

Proof. See Lemma 2.2 in ². □

The following result should be also known.

Theorem 8.1.2 (Continuous mapping). *Let $g : \mathcal{X} \rightarrow \mathbb{R}^m$ be continuous.*

(i) *If $X_n \rightsquigarrow X$ then $g(X_n) \rightsquigarrow g(X)$.*

(ii) *If $X_n \xrightarrow{p} X$ then $g(X_n) \xrightarrow{p} g(X)$.*

Proof. (i). Let $U \subseteq \mathbb{R}^m$ be an open set. Since g is continuous, $g^{-1}(U)$ is open in \mathcal{X} . By Portmanteau lemma,

$$\begin{aligned} \liminf \mathbb{P}(g(X_n) \in U) &= \liminf \mathbb{P}(X_n \in g^{-1}(U)) \\ &\geq \mathbb{P}(X \in g^{-1}(U)) = \mathbb{P}(g(X) \in U). \end{aligned}$$

Using the Portmanteau lemma again, we conclude that $g(X_n) \rightsquigarrow g(X)$.

(ii). Fix $\epsilon > 0$. For each $\delta > 0$ let B_δ be the set of x for which there exists y with $d(x, y) < \delta$, but $d(g(x), g(y)) > \epsilon$. If $X \notin B_\delta$ and $d(g(X_n), g(X)) > \epsilon$, then $d(X_n, X) \geq \delta$. Consequently,

$$\mathbb{P}(d(g(X_n), g(X)) > \epsilon) \leq \mathbb{P}(X \in B_\delta) + \mathbb{P}(d(X_n, X) \geq \delta).$$

The second term on the right converges to zero as $n \rightarrow \infty$ for every fixed $\delta > 0$. Because $B_\delta \cap \mathcal{X} \downarrow \emptyset$ by continuity of g , the first term converges to zero as $\delta \rightarrow 0$. □

It is important to remember basic relations between different notions of convergence.

Theorem 8.1.3. *Let X_n, X and Y_n be random vectors. Then*

(i) $X_n \xrightarrow{p} X$ implies $X_n \rightsquigarrow X$.

(ii) $X_n \xrightarrow{p} c$ for $c \in \mathbb{R}$ if and only if $X_n \rightsquigarrow c$.

(iii) if $X_n \rightsquigarrow X$ and $d(X_n, Y_n) \xrightarrow{p} 0$, then $Y_n \rightsquigarrow X$.

Moreover, if $X_n \rightsquigarrow X$ and $Y_n \rightarrow c$ then

(iv) $X_n + Y_n \rightsquigarrow X + c$,

(v) $X_n Y_n \rightsquigarrow cX$.

(vi) $Y_n^{-1} X_n \rightsquigarrow c^{-1} X$ provided $c \neq 0$.

The second part of the theorem is called the Slutsky lemma. For a proof see Theorem 2.7 and Lemma 2.8 in ³.

² A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998

³ A. W. van der Vaart. *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 1998

← Exercise ??

Example 8.1.4 (t-statistic). Say (X_n) is a series of i.i.d. random variables with $\mathbb{E}X_i = \mu$ and $\text{var}(X_i) = \sigma^2$. Let $\bar{X}_n = \frac{1}{n} \sum_i X_i$ and let

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

Then

$$t_{n-1} := \frac{\bar{X}_n - \mu}{S_n} \sqrt{n-1} \rightsquigarrow N(0, 1).$$

To prove this, note that by CLT we have $\sqrt{n}(\bar{X}_n - \mu) \rightsquigarrow N(0, \sigma^2)$. Also, by the law of large numbers $\bar{X}_n \xrightarrow{p} \mu$ and $\frac{1}{n} \sum X_i^2 \xrightarrow{p} \mathbb{E}X_i^2 = \sigma^2 + \mu^2$. By the continuous mapping theorem $\bar{X}_n^2 \xrightarrow{p} \mu^2$, which together with the Slutsky lemma implies that

$$S_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \rightsquigarrow (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

Hence, we have

$$\sqrt{\frac{n-1}{n}} \frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} \rightsquigarrow \frac{1}{\sigma} N(0, \sigma^2) = N(0, 1).$$

We say that X_n is **bounded in probability** if for every $\epsilon > 0$ there exists $M \in \mathbb{R}$ such that

$$\sup_n \mathbb{P}(\|X_n\| > M) < \epsilon.$$

By Prohorov’s theorem if $X_n \rightsquigarrow X$ for some X then it is bounded in probability (prove it!). Conversely, if X_n is bounded in probability then for some subsequence X_{n_k} we have $X_{n_k} \rightarrow X$ for some X .

We now introduce a special notation. Write $X_n = o_P(1)$ if $X_n \xrightarrow{p} 0$ and $X_n = O_P(1)$ if X_n is bounded in probability. More generally,

← Exercise 8.6.2

$$X_n = o_P(R_n) \quad \text{means} \quad X_n = Y_n R_n \text{ and } Y_n = o_P(1),$$

$$X_n = O_P(R_n) \quad \text{means} \quad X_n = Y_n R_n \text{ and } Y_n = O_P(1).$$

Remark 8.1.5. It is clear that $X_n = o_P(1)$ implies $X_n = O_P(1)$. Indeed, fix $\epsilon > 0$, if $X_n = o_P(1)$ then, for any $M_0 > 0$, there exist $N \in \mathbb{N}$ such that $\mathbb{P}(\|X_n\| > M_0) < \epsilon$ for all $n \neq N$. Let M_1, \dots, M_{N-1} be any numbers such that $\mathbb{P}(\|X_n\| > M_n) < \epsilon$. Taking $M := \max\{M_0, M_1, \dots, M_{N-1}\}$ we get that $\mathbb{P}(\|X_n\| > M) < \epsilon$ for all $n \in \mathbb{N}$.

← Exercise 8.6.3

← Exercise 8.6.4

The following result allows us to effectively work with Taylor series expansions.

Lemma 8.1.6. Let $R : \mathcal{X} \rightarrow \mathbb{R}$ be a continuous function and let X_n be a sequence of random vectors with values in \mathcal{X} and such that $X_n = o_P(1)$.

Then, for every $p > 0$,

(i) if $R(h) = o(\|h\|^p)$ as $h \rightarrow 0$, then $R(X_n) = o_P(\|X_n\|^p)$,

(ii) if $R(h) = O(\|h\|^p)$ as $h \rightarrow 0$, then $R(X_n) = O_P(\|X_n\|^p)$.

Proof. Define $g(h)$ as $g(h) = R(h)/\|h\|^p$ for $h \neq 0$ and $g(0) = 0$. Then $R(X_n) = g(X_n)\|X_n\|^p$.

(i) By assumption, g is continuous at zero, and so it is continuous everywhere. By the continuous mapping theorem, $g(X_n) \xrightarrow{P} g(0) = 0$.

(ii) By assumption, there exists M_0 and $\delta > 0$ such that $|g(h)| \leq M_0$ whenever $\|h\| \leq \delta$. Thus, for every n , $\mathbb{P}(|g(X_n)| > M_0) \leq \mathbb{P}(\|X_n\| > \delta) \rightarrow 0$, and so the sequence $g(X_n)$ is bounded in probability by exactly the same argument as we used in Remark 8.1.5. □

8.2 Delta Method

Suppose an estimator T_n for a parameter θ is available but the quantity of interest is $\phi(\theta)$ for some known function ϕ . It is natural to use the plug-in estimator $\phi(T_n)$. For example, if T_n is the MLE for θ then $\phi(T_n)$ is the MLE for $\phi(\theta)$. But how do the asymptotic properties of $\phi(T_n)$ follow those of T_n ?

The continuous-mapping theorem implies that if T_n is consistent for θ and ϕ is continuous, then $\phi(T_n)$ is consistent for $\phi(\theta)$. Here we show that a much stronger statement is true if ϕ is differentiable: if $\sqrt{n}(T_n - \theta) \rightsquigarrow T$ then $\sqrt{n}(\phi(T_n) - \phi(\theta)) \rightsquigarrow \phi'_\theta(T)$, where ϕ'_θ is the linear mapping representing the derivative of ϕ at θ (cf. Section A.3). In particular, asymptotic normality is preserved as a linear transformation of a Gaussian vector is Gaussian.

Theorem 8.2.1. *Let $U \subset \mathbb{R}^k$ open and let $\phi : U \rightarrow \mathbb{R}^m$ be a map differentiable at $\theta \in U$. Let T_n be random vectors taking values in U . If $r_n(T_n - \theta) \rightsquigarrow T$ for numbers $r_n \rightarrow \infty$, then $r_n(\phi(T_n) - \phi(\theta)) \rightsquigarrow \phi'_\theta(T)$. Moreover, the difference between $r_n(\phi(T_n) - \phi(\theta))$ and $\phi'_\theta(r_n(T_n - \theta))$ converges to zero in probability.*

Proof. Because $r_n(T_n - \theta)$ converges in distribution, we have $r_n(T_n - \theta) = O_P(1)$ and so $T_n - \theta = o_P(1)$. By differentiability of ϕ , $R(h) = \phi(\theta + h) - \phi(\theta) - \phi'(h)$ satisfies $R(h) = o(\|h\|)$ as $h \rightarrow \mathbf{0}$. Lemma 8.1.6 allows to replace the fixed h by a random sequence and gives

$$\phi(T_n) - \phi(\theta) - \phi'_\theta(T_n - \theta) \equiv R(T_n - \theta) = o_P(\|T_n - \theta\|).$$

Multiply this with r_n and note that $o_P(r_n\|T_n - \theta\|) = o_P(1)$ because $r_n(T_n - \theta) = O_P(1)$. This gives the last statement of the theorem. Because linear maps are continuous, the continuous mapping theorem

gives that $\phi'_\theta(r_n(T_n - \theta)) \rightsquigarrow \phi'_\theta(T)$. By the Theorem 8.1.3(iii) we get that $r_n(\phi(T_n) - \phi(\theta))$ has the same weak limit. \square

Some basic examples of that you saw in the first semester.

Example 8.2.2. *The delta method may be useful in a wide range of scenarios. For example, suppose that we want to obtain asymptotic confidence intervals for some parameter but the asymptotic distribution depends on θ . More concretely, suppose X_1, \dots, X_n are i.i.d. $\text{Pois}(\theta)$. Then $\sqrt{n}(\bar{X}_n - \theta) \rightsquigarrow N(0, \theta)$. By the delta method $\sqrt{n}(f(\bar{X}_n) - f(\theta)) \rightsquigarrow N(0, (f'(\theta))^2 \theta)$. If we solve for $f'(\theta) = C/\sqrt{\theta}$ we get $f(\theta) = \sqrt{\theta}$, $C = 1/2$ and then*

$$\sqrt{n}(\sqrt{\bar{X}_n} - \sqrt{\theta}) \rightsquigarrow N(0, \frac{1}{4}).$$

Therefore we can easily construct the asymptotic confidence interval for $\sqrt{\theta}$ as $\sqrt{\bar{X}_n} \pm \frac{z_{\alpha/2}}{4\sqrt{\bar{X}_n}}$. This example easily generalizes.

← Exercise 8.6.5

As an important application of the delta method we show that the MLE in the exponential family has an asymptotically Gaussian distribution. Recall that in the exponential family (1.1) it holds that

$$\mathbb{E}_\theta(\mathbf{t}(X)) = \mu(\theta) = \nabla A(\theta), \quad \text{var}_\theta(\mathbf{t}(X)) = V(\theta) = \nabla^2 A(\theta).$$

Suppose that the random sample X_1, \dots, X_n is generated from \mathbb{P}_{θ_0} . Denote $\bar{\mathbf{t}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{t}(X_i)$. By the central limit theorem, it follows that

$$\sqrt{n} \nabla \ell_n(\theta_0) = \sqrt{n}(\bar{\mathbf{t}}_n - \mu(\theta_0)) \rightsquigarrow N_d(\mathbf{0}, V(\theta_0)). \quad (8.1)$$

Using the delta method we conclude.

Theorem 8.2.3. *The MLE $\hat{\theta}_n$ in the exponential family (1.1) based on the sample X_1, \dots, X_n satisfies $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(\mathbf{0}, V^{-1}(\theta_0))$.*

Proof. As noted above $\sqrt{n}(\bar{\mathbf{t}}_n - \mu(\theta_0)) \rightsquigarrow N(\mathbf{0}, V(\theta_0))$. The delta method, applied with $\phi = \mu^{-1}$, implies that $\sqrt{n}(\hat{\theta}_n - \theta_0)$ also converges in distribution to a Gaussian distribution. The fact that $\hat{\theta}_n = \mu^{-1}(\bar{\mathbf{t}}_n)$ is part of Proposition 1.3.6. The Jacobian of μ^{-1} at $\mu(\theta_0)$ is the inverse of the Jacobian of μ at θ_0 , which is equal to $\nabla^2 A(\theta_0) = V(\theta_0)$. Thus, $\phi'_{\mu(\theta_0)}(T) = V^{-1}(\theta_0) \cdot T$. Now it is straightforward to check that the asymptotic covariance matrix is $V^{-1}(\theta_0)$. \square

With small amount of extra work, this result can be generalized to curved exponential families for which the canonical parameter θ is parametrized in a smooth way by a lower dimensional parameters τ . Denote by $\hat{\tau}$ the MLE in this smaller model and assume that the sample is generated from the parameter $\theta_0 = \theta(\tau_0)$ corresponding to τ_0 . Taking the first-order expansion around τ_0 gives

← Exercise 8.6.6

$$0 = \nabla_{\boldsymbol{\tau}} \ell_n(\boldsymbol{\theta}(\widehat{\boldsymbol{\tau}}_n)) = \nabla_{\boldsymbol{\tau}} \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0)) + \nabla_{\boldsymbol{\tau}}^2 \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0))(\widehat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) + o_P(\|\boldsymbol{\tau}_0 - \widehat{\boldsymbol{\tau}}_n\|).$$

Multiply now this equation by \sqrt{n} and use Exercise 8.6.6 to conclude that

$$-\nabla_{\boldsymbol{\tau}}^2 \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0)) \cdot \sqrt{n}(\widehat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0) = \sqrt{n} \nabla_{\boldsymbol{\tau}} \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0)) + o_P(1). \quad (8.2)$$

By the chain rule, $\nabla_{\boldsymbol{\tau}} \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0)) = (\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\tau}}(\boldsymbol{\tau}_0))^T \nabla \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0))$, and so, using (8.1), we get

$$\sqrt{n} \nabla_{\boldsymbol{\tau}} \ell_n(\boldsymbol{\theta}(\boldsymbol{\tau}_0)) \rightsquigarrow N\left(\mathbf{0}, \left(\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\tau}}(\boldsymbol{\tau}_0)\right)^T \cdot V(\boldsymbol{\theta}_0) \cdot \frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\tau}}(\boldsymbol{\tau}_0)\right).$$

By Slutsky Lemma (Theorem 8.1.3(iv)), the left-side expression in (8.2) also converges in distribution to the same Gaussian. By Exercise 8.6.7 we conclude that $\sqrt{n}(\widehat{\boldsymbol{\tau}}_n - \boldsymbol{\tau}_0)$ is asymptotically normal with mean zero and the covariance matrix equal to the *inverse* of the Fisher information matrix $\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\tau}}(\boldsymbol{\tau}_0)^T \cdot V(\boldsymbol{\theta}_0) \cdot \frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\tau}}(\boldsymbol{\tau}_0)$.

← Exercise 8.6.7

8.3 M-estimators

In this chapter X_1, \dots, X_n are i.i.d. \mathbb{P} where $\mathbb{P} \in \mathcal{P}$, with $\mathcal{P} = \{\mathbb{P}_\theta : \theta \in \Theta\}$ such that

1. $\Theta \subset \mathbb{R}^d$ open,
2. densities of \mathbb{P}_θ are $p(\cdot, \theta)$, $\theta \in \Theta$.

M-estimators were defined and discussed in Section 4.1.2. In this section we discuss basic asymptotics of this class of estimators.

8.3.1 Consistency

In order to develop some general asymptotic theory of M-estimators we start by discussing their consistency. Under minor conditions the law of large numbers will give us that

$$M_n(\theta) \xrightarrow{P} M(\theta) \quad \text{for every } \theta, \quad (8.3)$$

where $M(\theta) = \mathbb{E}m_\theta(X)$. It is reasonable to expect that the sequence of maximizers $\hat{\theta}_n$ of M_n converges in probability to the maximizer of M . However, pointwise convergence in (8.3) is too weak because $\hat{\theta}_n$ depends on the whole function $\theta \mapsto M_n(\theta)$. We present approach based on the assumption of the **uniform convergence** of M_n to M : $\|M_n - M\|_\infty \xrightarrow{P} 0$, where

$$\|M_n - M\|_\infty := \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)|.$$

Theorem 8.3.1. Let $\{M_n\}$ be a sequence of random functions, continuous on Θ . Say $\|M_n - M\|_\infty \xrightarrow{P} 0$, where M is some non-random continuous function on Θ . Then

1. If $T_n \xrightarrow{P} T^*$ then $M_n(T_n) \xrightarrow{P} M(T^*)$.
2. If $t^* = \arg \max M(t)$ uniquely and $T_n \in \arg \max M_n(t)$, then $T_n \xrightarrow{P} t^*$.

Proof. 1. By the triangle inequality

$$|M_n(T_n) - M(T^*)| \leq |M_n(T_n) - M(T_n)| + |M(T_n) - M(T^*)|.$$

The second term on the right goes to zero in probability by the continuous mapping theorem. The first term goes to zero by the uniform convergence hypothesis in the theorem because $|M_n(T_n) - M(T_n)| \leq \|M_n - M\|_\infty \xrightarrow{P} 0$. We conclude that $|M_n(T_n) - M(T^*)| \xrightarrow{P} 0$ or, equivalently $M_n(T_n) \xrightarrow{P} M(T^*)$. 2. For some $\epsilon > 0$, let $K_\epsilon = K \cap (B_\epsilon(t^*))^c$, where $B_\epsilon(t^*)$ is a ball of radius ϵ around t^* . Let $m = M(t^*)$ and let $m_\epsilon = \sup_{t \in K_\epsilon} M(t)$. Since t^* is unique $\delta := m - m_\epsilon > 0$. From uniform convergence, there is some N s.t. for all $n > N$ we have $\|M_n - M\|_\infty < \frac{\delta}{2}$. We can therefore infer the following:

$$\sup_{t \in K_\epsilon} M_n(t) < \sup_{t \in K_\epsilon} M(t) + \frac{\delta}{2} = m_\epsilon + \frac{\delta}{2} = m - \frac{\delta}{2}$$

and

$$M_n(T_n) = \max_{t \in K} M_n(t) \geq M_n(t^*) > M(t^*) - \frac{\delta}{2} = m - \frac{\delta}{2}.$$

By these two inequalities we get that $T_n \notin K_\epsilon$, and thus $T_n \in B_\epsilon(t^*)$. Hence, $\|M_n - M\|_\infty < \delta \Rightarrow \|T_n - t^*\| < \epsilon \Rightarrow \mathbb{P}(\|T_n - t^*\| > \epsilon) \leq \mathbb{P}(\|M_n - M\|_\infty > \frac{\delta}{2}) \Rightarrow T_n \xrightarrow{P} t^*$. \square

The assumption of uniform convergence holds for the log-likelihood function in exponential families but is typically too strong, especially when Θ is not compact.

Example 8.3.2. Consider the problem of estimating the variance using the sample variance $s_n^2 = \frac{1}{n} \sum_i x_i^2$. This corresponds to optimizing the function

$$M_n(\sigma) = \log \sigma + \frac{s_n^2}{2\sigma^2}$$

with $M(\sigma) = \log \sigma + \frac{\sigma_0^2}{2\sigma^2}$, where σ_0^2 is the variance of the sample. Note that $|M_n(\sigma) - M(\sigma)|$ is unbounded for $\sigma > 0$. Nevertheless, consistency follows directly by the law of large numbers.

Under some regularity conditions it is possible to get a fairly universal approach that does not require uniform convergence. Suppose that $m_\theta(x)$ is twice continuously differentiable with respect to θ on the open parameters set Θ and consider a Taylor series expansion of ∇M_n at the true data generating θ_0 . Since Θ is open then the M-estimator $\hat{\theta}_n$ satisfies $\nabla M_n(\hat{\theta}_n) = \mathbf{0}$ and so

$$\mathbf{0} = \nabla M_n(\theta_0) + \nabla^2 M_n(\tilde{\theta}_n)(\hat{\theta}_n - \theta_0), \quad (8.4)$$

where $\tilde{\theta}_n$ lies between $\hat{\theta}_n$ and θ_0 . By the law of large numbers, as long as $\mathbb{E}_{\theta_0} \nabla m_{\theta_0}(X) < \infty$,

$$\nabla M_n(\theta_0) \xrightarrow{P} \mathbb{E}_{\theta_0} \nabla M_n(\theta_0).$$

To establish consistency of $\hat{\theta}_n$ using (8.4) it suffices to have that:

← Exercise 8.6.8

- (i) $\mathbb{E}_{\theta_0} \nabla M_n(\theta_0) = 0$ and
- (ii) $(\nabla^2 M_n(\tilde{\theta}_n))^{-1} = O_P(1)$.

In Exercise 8.6.8 we show that (i) holds when $m_\theta(x) = -\log p_\theta(x)$ (MLE) as well as $m_\theta(x) = \|\theta - \delta(X)\|^2$, where $\delta(X)$ is an unbiased estimator of the parameter θ . More generally, (i) holds if θ_0 minimizes the function $\theta \mapsto \mathbb{E}_{\theta_0} m_\theta(X)$. Thus, typically (ii) remains the main regularity assumption that needs to be checked in order to establish consistency.

8.3.2 Asymptotic normality

We will now discuss asymptotic normality of M-estimators assuming that $\hat{\theta}_n$ is a consistent estimator and $\Theta \subseteq \mathbb{R}^d$ is open. If ψ_θ in (4.5) is twice differentiable then there is a standard approach for proving asymptotic normality based on the Taylor series expansion. As earlier, suppose that

$$\Psi_n(\theta) = \frac{1}{n} \sum_{i=1}^n \psi_\theta(X_i), \quad \Psi(\theta) = \mathbb{E} \psi_\theta(X) = \mathbb{E} \Psi_n(\theta)$$

and assume that $\hat{\theta} \xrightarrow{P} \theta_0$, where

$$\Psi_n(\hat{\theta}_n) = 0, \quad \Psi(\theta_0) = 0.$$

Taylor's theorem gives

$$0 = \Psi_n(\hat{\theta}_n) = \Psi_n(\theta_0) + \nabla \Psi_n(\theta_0)(\hat{\theta}_n - \theta_0) + R_n,$$

where $R_n = o_P(\|\hat{\theta}_n - \theta_0\|)$. In other words,

$$\nabla \Psi_n(\theta_0) \sqrt{n}(\hat{\theta}_n - \theta_0) = -\sqrt{n} \Psi_n(\theta_0) - \sqrt{n} R_n,$$

By the central limit theorem, as long as, $\mathbf{B}(\theta_0) := \mathbb{E}\Psi_{\theta_0}(X)\Psi_{\theta_0}^\top(X)$ has finite entries,

$$-\sqrt{n}\Psi_n(\hat{\theta}_n) \rightsquigarrow N(0, \mathbf{B}(\theta_0)).$$

By the law of large numbers, as long as $\mathbf{A}(\theta_0) := \mathbb{E}\nabla\psi_{\theta_0}(X)$ has finite entries, then

$$\nabla\Psi_n(\theta_0) \xrightarrow{p} \mathbf{A}(\theta_0)$$

Under suitable regularity conditions $\sqrt{n}R_n = o_P(1)$, in which case we can use the Slutsky lemma to conclude that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, \mathbf{A}(\theta_0)^{-1}\mathbf{B}(\theta_0)\mathbf{A}(\theta_0)^{-\top}). \quad (8.5)$$

The conditions that assure that $\sqrt{n}R_n = o_P(1)$ may be complicated in general. However, they will hold if Ψ_n is sufficiently smooth.

A special case of interest is the maximum likelihood estimator. In this case $\psi_\theta(x) = \nabla \log p_\theta(x)$ and so

$$\mathbb{E}\psi_\theta(X) = \mathbb{E}\nabla\ell_n(\theta), \quad \mathbb{E}\nabla\psi_\theta(X) = \mathbb{E}\nabla^2\ell_n(\theta).$$

Under some regularity conditions the maximum likelihood estimator is consistent and asymptotically normal: as in (8.5), with

$$\mathbf{A}(\theta_0) = -\mathbb{E}_{\theta_0}[\nabla^2\ell_n(\theta_0)] \quad \text{and} \quad \mathbf{B}(\theta_0) = \mathbb{E}_{\theta_0}[\nabla\ell_n(\theta_0)\nabla\ell_n(\theta_0)^\top].$$

In this case $\mathbf{B}(\theta_0)$ is the variance of the score also known as the Fisher information matrix. If $\log p_\theta(x)$ is twice continuously differentiable, we can write

$$\nabla^2 \log p_\theta(X) = \frac{1}{p_\theta(X)} \nabla^2 p_\theta(X) - (\nabla \log p_\theta(X))(\nabla \log p_\theta(X))^\top.$$

Since

$$\mathbb{E}_\theta \frac{1}{p_\theta(X)} \nabla^2 p_\theta(X) = \int \nabla^2 p_\theta(x) dx = \nabla^2 \int p_\theta(x) dx = 0,$$

we conclude that in this case $\mathbf{A}(\theta) = \mathbf{B}(\theta)$ and so the MLE $\hat{\theta}_n$ satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, \mathbf{B}(\theta_0)^{-1}). \quad (8.6)$$

It is of considerable interest to establish asymptotic normality also in the case when Ψ_n is not twice differentiable. But we will not discuss this case in more detail here.

← Exercise 8.6.9

8.4 Generalized likelihood ratio test

The tests in Chapter 3 have strong optimality properties but require conditions on the densities for the data and the form of the hypotheses that are rather special and can fail for many natural models.

By contrast, the generalized likelihood ratio test introduced in this chapter requires little structure, but it does not have exact optimality properties. Use of this test is justified by large sample theory.

Let the data X_1, \dots, X_n be i.i.d. with common density p_θ for $\theta \in \Theta$. We want to test $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$, where $\Theta_0 \subset \Theta_1 \subseteq \Theta$. Note that the hypotheses are nested.

A sensible extension of the idea behind the likelihood ratio test, as discussed in Proposition 3.1.1, is to base a test on the log-likelihood ratio

$$\lambda_n = \lambda_n(X_1, \dots, X_n) = \log \frac{\sup_{\theta \in \Theta_1} \prod_{i=1}^n p_\theta(X_i)}{\sup_{\theta \in \Theta_0} \prod_{i=1}^n p_\theta(X_i)}.$$

As before, the null hypothesis is rejected for large values of the statistic. Let $\ell_n(\theta) = \frac{1}{n} \sum_i \log p_\theta(X_i)$ be the log-likelihood function and suppose both suprema are attained. If $\hat{\theta}_n$ is the optimum over Θ_1 and $\hat{\theta}_{n,0}$ is the optimum over Θ_0 then

$$\lambda_n = n(\ell_n(\hat{\theta}_n) - \ell_n(\hat{\theta}_{n,0})).$$

We next give an example that can be viewed as the limiting situation for which the approximation is exact:

Example 8.4.1 (The Gaussian sequence model). Let Y_1, \dots, Y_n be independent with $Y_i \sim N(\mu_i, \sigma_0^2)$ where σ_0 is known. In other words,

$$\mathbf{Y} = (Y_1, \dots, Y_n) \sim N_n(\boldsymbol{\mu}, \sigma_0^2 I_n).$$

This model is sometimes called the Gaussian sequence model and will be one of the important examples in our discussion of high-dimensional problems. We are now interested in testing whether $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$ is a member of a q -dimensional linear subspace $\mathcal{L}_0 \subset \mathbb{R}^n$, versus the alternative that $\boldsymbol{\mu} \in \mathcal{L} \setminus \mathcal{L}_0$ where \mathcal{L} is an r -dimensional linear subspace of \mathbb{R}^d and $\mathcal{L}_0 \subset \mathcal{L}$, $r > q$.

Transform to canonical form by setting $\mathbf{U} = \mathbf{Q}\mathbf{Y}$, where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix with the first q rows spanning \mathcal{L}_0 and the first r rows spanning \mathcal{L} . By construction, $\mathbf{U} \sim N_n(\mathbf{Q}\boldsymbol{\mu}, \sigma_0^2 I_n)$, where the mean vector $\boldsymbol{\eta} = \mathbf{Q}\boldsymbol{\mu}$ satisfies $\eta_i = 0$ for $i \geq q+1, \dots, n$ under H_0 , and $\eta_i = 0$ for $i \geq r+1, \dots, n$ under H_1 .

Set $\theta_i = \eta_i/\sigma_0$, $i = 1, \dots, r$ and $X_i = U_i/\sigma_0$, $i = 1, \dots, n$. Then $\mathbf{X} \sim N(\boldsymbol{\theta}, I_n)$. Moreover, the hypothesis H_0 is equivalent to $H_0 : \theta_i = 0$ for $i \geq q+1$, and the alternative is $H_1 : \theta_i = 0$ for $i \geq r+1$. The log-likelihood for $\mathbf{X} \sim N(\boldsymbol{\theta}, I_n)$ is

$$\ell(\boldsymbol{\theta}) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\theta}\|^2,$$

and so $\arg \max_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta}; \mathbf{x}) = \arg \min_{\boldsymbol{\theta} \in \Theta} \|\mathbf{x} - \boldsymbol{\theta}\|^2$. Under H_0 ,

$$2\lambda = 2(\ell(\hat{\boldsymbol{\theta}}) - \ell(\hat{\boldsymbol{\theta}}_0)) = \sum_{i=q+1}^r X_i^2 \sim \chi_{r-q}^2.$$

It is a remarkable fact that χ_{r-1}^2 holds as an approximation to the null distribution of $2\lambda_n$ quite generally when the hypothesis is a q -dimensional submanifold of an r -dimensional parameter space. Some of the sufficient regularity conditions under which this holds were briefly discussed in Section 8.3.1 (e.g. Θ open, $\log p_\theta(x)$ twice continuously differentiable with respect to θ , some conditions on the behavior of the hessian $\nabla^2 \ell_n$ around $\theta_0 \in \Theta_0$).

Theorem 8.4.2 (Wilk's theorem). *Consider testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$, where $\Theta_0 \subset \Theta_1 \subseteq \Theta$ and Θ_0, Θ_1 are manifolds of dimension q and r respectively. Under the same conditions guaranteeing asymptotic normality of the two MLEs, we have under the null hypothesis,*

$$2\lambda_n \rightsquigarrow \chi_{r-q}^2.$$

Sketch of the proof. For simplicity, we assume that $\Theta_1 = \Theta$. Let $\hat{\theta}_{0,n}$ be the MLE over Θ_0 and $\hat{\theta}_n$ be the MLE over Θ . Directly by definition

$$2\lambda_n = 2n(\ell_n(\hat{\theta}_n) - \ell_n(\hat{\theta}_{0,n})).$$

Applying the first order Taylor expansion of ℓ_n at $\hat{\theta}_n$, and using the fact that $\nabla \ell_n(\hat{\theta}_n) = 0$, we find that

$$\ell_n(\hat{\theta}_n) - \ell_n(\hat{\theta}_{0,n}) = -\frac{1}{2}(\hat{\theta}_n - \hat{\theta}_{0,n})^\top \nabla_{\tilde{\theta}}^2 \ell_n(\tilde{\theta}_n)(\hat{\theta}_n - \hat{\theta}_{0,n}),$$

where $\tilde{\theta}_n$ lies between $\hat{\theta}_{0,n}$ and $\hat{\theta}_n$. By asymptotic normality $\hat{\theta}_{0,n} \xrightarrow{p} \theta_0$ and $\hat{\theta}_n \xrightarrow{p} \theta_0$ and so also $\tilde{\theta}_n \xrightarrow{p} \theta_0$. Hence, we argue that

$$-\nabla_{\tilde{\theta}}^2 \ell_n(\tilde{\theta}_n) \xrightarrow{p} -\mathbb{E}[\nabla_{\tilde{\theta}}^2 \ell_n(\theta_0)] = I(\theta_0).$$

(This statement is actually a bit subtle and may require some further regularity assumptions.) Therefore

$$2n(\ell_n(\hat{\theta}_n) - \ell_n(\hat{\theta}_{0,n})) = \sqrt{n}(\hat{\theta}_n - \hat{\theta}_{0,n})^\top \left(-\nabla_{\tilde{\theta}}^2 \ell_n(\tilde{\theta}_n) \right) \sqrt{n}(\hat{\theta}_n - \hat{\theta}_{0,n}).$$

In the case when $r = d$, $\hat{\theta}_{0,n} = \theta_0$ since there is no maximization and then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow N(0, I(\theta_0)^{-1})$$

from the asymptotic normality of the MLE. Now putting together the pieces:

$$\sqrt{n}(\hat{\theta}_n - \theta_0)^\top \left(-\frac{1}{n} \nabla_{\tilde{\theta}}^2 \ell_n(\tilde{\theta}_n) \right) \sqrt{n}(\hat{\theta}_n - \theta_0) \rightsquigarrow \sum_{i=1}^d Z_i^2 = \chi_d^2,$$

where $Z_i \sim N(0, 1)$. The case when $r < d$ does not give much more insight, but a lot more algebra. \square

Example 8.4.3 (Exponential families). Suppose that the observations are sampled from a density p_θ in the d -dimensional regular exponential family

$$p_\theta(x) = h(x)e^{\langle \theta, t(x) \rangle - A(\theta)}.$$

Let $\Theta \subseteq \mathbb{R}^d$ be the canonical parameter space, and consider testing a null hypothesis $\Theta_0 \subset \Theta$ versus its complement. The log-likelihood ratio statistic is given by

$$\lambda_n = n \sup_{\theta \in \Theta} \inf_{\theta \in \Theta_0} (\langle \theta - \theta_0, \bar{t}_n \rangle - A(\theta) + A(\theta_0)).$$

Note that the Kullback-Leibler divergence of the measures \mathbb{P}_{θ_0} and \mathbb{P}_θ is equal to

$$K(\theta, \theta_0) = \mathbb{E}_\theta \log \frac{p_\theta}{p_{\theta_0}} = \langle \theta - \theta_0, \mathbb{E}_\theta t(X) \rangle - A(\theta) + A(\theta_0).$$

If the maximum likelihood estimator $\hat{\theta} \in \Theta$ exists, then $\mathbb{E}_{\hat{\theta}}(t(X)) = \bar{t}_n$ by Proposition 1.3.6. Comparing the two preceding displays, we see that the likelihood ratio statistic can be written as

$$\lambda_n = n \inf_{\theta_0 \in \Theta_0} K(\hat{\theta}, \theta_0).$$

This formula can be used to study the asymptotic properties of the likelihood ratio statistic directly.

8.5 Limits of Bayesian procedures

In this section θ will denote a random parameter rather than a vector. If θ possesses a density π , then the density of the posterior distribution of θ is given by Bayes' formula

$$\pi(\theta | X_1, \dots, X_n) = \frac{\prod_{i=1}^n p(X_i | \theta) \pi(\theta)}{\int \prod_{i=1}^n p(X_i | \theta) \pi(\theta) d\theta} \quad (8.7)$$

which is a random kernel in the same sense as randomized decision rules; c.f. Definition 2.2.1. This expression may define a probability density even if π is not a probability density itself. A prior distribution with infinite mass is called improper.

We pose the question, under what conditions are Bayes methods asymptotically consistent/efficient. We start with an example.

Example 8.5.1. $X \sim \text{Bin}(n, \theta)$ with prior $\theta \sim \text{Beta}(a, b)$, that is, $\pi(\theta) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \theta^{a-1} (1-\theta)^{b-1}$. Since the posterior satisfies

$$\pi(\theta | x) \propto \theta^{a+x-1} (1-\theta)^{b+n-x-1},$$

it must be equal to $\text{Beta}(a+x, b+n-x)$. Since the mean in $\text{Beta}(a, b)$ is $\frac{a}{a+b}$, the Bayes estimate under the squared loss x (c.f. Proposition 2.3.5) is:

$$\delta_{a,b}^n = \mathbb{E}(\theta | x) = \frac{a+x}{a+b+n}.$$

If we assume that the true value of θ is θ_0 ,

$$\sqrt{n}(\delta_{a,b}^n - \theta_0) = \sqrt{n} \left(\frac{X}{n} - \theta_0 \right) + \frac{\sqrt{n}}{a+b+n} \left(a - (a+b) \frac{X}{n} \right).$$

Note that $\sqrt{n}(\frac{X}{n} - \theta_0) \rightsquigarrow N(0, \theta_0(1 - \theta_0))$ since $\frac{X}{n}$ is MLE of θ and

$$\frac{\sqrt{n}}{a+b+n} \left(a - (a+b) \frac{X}{n} \right) \xrightarrow{p} 0.$$

Therefore the Bayes estimate is asymptotically consistent and efficient for all a and b .

Bernstein-von Mises theorem guarantees that the random kernel $\pi(\theta|X_1, \dots, X_n)$ in (8.7) is “close” to $N(\hat{\theta}_n, (nI(\theta_0))^{-1})$, where $\hat{\theta}_n$ is the MLE and θ_0 is the true value.

Theorem 8.5.2 (Bernstein-von Mises). *Assume that the prior density π is continuous and strictly positive and the standard regularity conditions for asymptotic normality of the MLE $\hat{\theta}_n$ hold. Then the conditional density of $\sqrt{n}(\theta - \hat{\theta}_n)$ given X_1, \dots, X_n converges to the PDF of $N(0, (I(\theta_0))^{-1})$.*

Proof. (Proof sketch, see Section 5.5 in ⁴ for details) Condition on X_1, \dots, X_n . Let $v = \sqrt{n}(\theta - \hat{\theta}_n)$ so that $\theta = \hat{\theta}_n + v/\sqrt{n}$. Denote

$$f(v) = \pi(\hat{\theta}_n + \frac{v}{\sqrt{n}} | X_1, \dots, X_n)$$

so that

$$\log f(v) = \text{const} + n\ell_n(\hat{\theta}_n + \frac{v}{\sqrt{n}}) + \log \pi(\hat{\theta}_n + \frac{v}{\sqrt{n}}). \quad (8.8)$$

Using the second order expansion of $\ell_n(\hat{\theta}_n + \frac{v}{\sqrt{n}})$ around $\hat{\theta}_n$ and using the fact that $\nabla \ell_n(\hat{\theta}_n) = 0$ we get

$$n\ell_n(\hat{\theta}_n + \frac{v}{\sqrt{n}}) = n\ell_n(\hat{\theta}_n) + \frac{1}{2}v^T \nabla^2 \ell_n(\tilde{\theta}_n)v,$$

where $\tilde{\theta}_n = \hat{\theta}_n + \frac{tv}{\sqrt{n}}$ for $t \in (0, 1)$. Plugging this to (8.8) gives

$$\log f(v) = \text{const} + \frac{1}{2}v^T \nabla^2 \ell_n(\hat{\theta}_n + \frac{tv}{\sqrt{n}})v + \log \pi(\hat{\theta}_n + \frac{v}{\sqrt{n}}). \quad (8.9)$$

This holds for any fixed $v \in \mathbb{R}^d$ and so, by (8.9), the conditional distribution of $v = \sqrt{n}(\theta - \hat{\theta}_n)$ is proportional to

$$\pi(\hat{\theta}_n + \frac{v}{\sqrt{n}}) \exp\left\{-\frac{1}{2}v^T (-\nabla^2 \ell_n(\hat{\theta}_n + \frac{tv}{\sqrt{n}}))v\right\}.$$

Since $\hat{\theta}_n$ converges in probability to θ_0 , v/\sqrt{n} converges to zero and under appropriate uniformity conditions, $-\nabla^2 \ell_n(\hat{\theta}_n + tv/\sqrt{n})$ converges to $I(\theta_0)$. Therefore, expression (8.9) after normalizing converges to the density of $N(0, (I(\theta_0))^{-1})$. \square

⁴Peter J. Bickel and Kjell A. Doksum. *Mathematical statistics—basic ideas and selected topics. Vol. 1.* Texts in Statistical Science Series. CRC Press, Boca Raton, FL, second edition, 2015

8.6 Exercises

Exercise 8.6.1. Say $X_n \rightsquigarrow X \in \mathbb{R}^d$ and $Y_n \rightsquigarrow c \in \mathbb{R}^k$ (constant) and let $f : \mathbb{R}^d \times \mathbb{R}^k \rightarrow \mathbb{R}^m$ continuous almost everywhere. Show that $f(X_n, Y_n) \rightsquigarrow f(X, c)$.

Exercise 8.6.2. Argue that $X_n \xrightarrow{P} X$ if and only if $d(X_n, X) = o_P(1)$.

Exercise 8.6.3. Show that:

$$\begin{aligned} o_P(1) + o_P(1) &= o_P(1) \\ o_P(1) + O_P(1) &= O_P(1) \\ o_P(1)O_P(1) &= o_P(1) \\ (1 + o_P(1))^{-1} &= O_P(1) \\ o_P(R_n) &= R_n o_P(1) \\ O_P(R_n) &= R_n O_P(1) \\ o_P(O_P(1)) &= o_P(1). \end{aligned}$$

Exercise 8.6.4. Let (X_n) be a random sequence in \mathbb{R}^m . Show that $X_n = o_P(1)$ if and only if each coordinate sequence is $o_P(1)$.

Exercise 8.6.5. Let X_1, \dots, X_n be i.i.d. $\mathbb{E}X_i = \mu$, $\text{var}(X_i) = 1$. Find constants such that $r_n(\bar{X}_n^2 - a_n)$ converges in distribution when $\mu = 0$ and when $\mu \neq 0$.

Exercise 8.6.6. Show that $\sqrt{n}(\hat{\tau}_n - \tau_0) = O_P(1)$. Conclude that $\hat{\tau}_n \xrightarrow{P} \tau_0$.

Exercise 8.6.7. Show that in exponential families

$$-\nabla_{\tau}^2 \ell_n(\theta(\tau_0)) \xrightarrow{P} \frac{\partial \theta}{\partial \tau}(\tau_0)^T \cdot V(\theta_0) \cdot \frac{\partial \theta}{\partial \tau}(\tau_0).$$

Exercise 8.6.8. Show that (i) in the bottom of Section 8.3.1 holds when $m_{\theta}(x) = -\log p_{\theta}(x)$ (MLE) as well as $m_{\theta}(x) = \|\theta - \delta(X)\|^2$, where $\delta(X)$ is an unbiased estimator of the parameter θ .

Exercise 8.6.9. Let X_1, \dots, X_n be a random sample from $N(\mu, 1)$, where it is known that $\mu \geq 0$. Show that the MLE is not asymptotically normal under $\mu = 0$. Why does this not contradict our result on asymptotic normality of the MLE?

Part IV

Mathematical appendix

A

Real Analysis

The Advanced Theory of Statistics is a technical subject. It is then important to actively look for high-level insights. This allows not only to understand the material better but also to see how the presented results may be generalized. Part of the goal of this lecture is to help students look for such insights.

In this appendix we briefly cover some fundamental concepts that help to understand many parts of theoretical statistics. These are: differentiation in vector spaces and convexity. This exposition assumes certain level of mathematical maturity on the level of a basic real analysis course. For more details, we refer to

<https://pzwiernik.github.io/docs/RealAnalysisNotes.pdf>.

A.1 Vector spaces

A set V is a **vector space** if (i) $\mathbf{0}$ lies in V (ii) for any two $\mathbf{x}, \mathbf{y} \in V$ also $\mathbf{x} + \mathbf{y} \in V$, (iii) if $\mathbf{x} \in V$ and $\lambda \in \mathbb{R}$ then $\lambda \cdot \mathbf{x}$ lies in V . A general abstract definition of a vector space is more complicated because it needs to explain what we mean by $\mathbf{0}$ and what we mean by the algebraic operations $+$ and \cdot . Here however, we always work with variations of the following three examples (Examples A.1.1-A.1.4), where all these objects are naturally defined.

Example A.1.1 (The Euclidean space \mathbb{R}^d). *The real space \mathbb{R}^d with elements $\mathbf{x} = (x_1, \dots, x_d)$, $x_i \in \mathbb{R}$, is an example of a vector space. We define the standard inner product as*

$$\langle \mathbf{x}, \mathbf{y} \rangle := x_1 y_1 + \dots + x_d y_d \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

The vector space \mathbb{R}^d equipped with the induced norm $\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ is called the Euclidean space. Every subset of \mathbb{R}^d given by linear equations also forms a vector space with the induced norm.

In general, an inner product on a vector space V is a function from

$V \times V$ to \mathbb{R} that must satisfy the following three conditions for all $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$

1. $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ (symmetry)
2. $\langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle = \alpha \langle \mathbf{x}, \mathbf{z} \rangle + \beta \langle \mathbf{y}, \mathbf{z} \rangle$ (linearity in the first argument)
3. $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ and is zero only if $\mathbf{x} = \mathbf{0}$. (positive definiteness)

The inner product space induces a norm in the standard way

$$\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

but not all norms are obtained in this way. Every vector space with a given norm becomes a metric space with distance function satisfying

$$d(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\| \quad \text{for } \mathbf{x}, \mathbf{y} \in V.$$

Example A.1.2 (The space of $m \times n$ matrices). The space $\mathbb{R}^{m \times n}$ of all real $m \times n$ matrices also forms a vector space. The standard inner product is

$$\langle A, B \rangle = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij} = \text{tr}(AB^T) \quad \text{for all } A, B \in \mathbb{R}^{m \times n}.$$

The induced norm is the Frobenius norm $\|A\|_F := \sqrt{\langle A, A \rangle}$ but for matrices we typically work with the operator norm instead

$$\|A\| := \max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\|.$$

A special case is the vector space \mathbb{S}^m of all $m \times m$ symmetric matrices with the inner product induced from $\mathbb{R}^{m \times m}$:

$$\mathbb{S}^m = \{A \in \mathbb{R}^{m \times m} : A = A^T\}.$$

Exercise A.1.3. Show that $\text{tr}(AB^T) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij}$ for all $A, B \in \mathbb{R}^{m \times n}$.

All finite dimensional vector spaces are similar and behave like the Euclidean space \mathbb{R}^d . In particular, defining a metric space structure is straightforward. In more complicated situations we may need to work a bit extra. Consider the following important example.

Example A.1.4 (L^2 functions). The set of all measurable functions $f : \mathcal{X} \rightarrow \mathbb{R}$ on some set \mathcal{X} also forms a vector space with addition defined pointwise (why? what is the zero of this vector space?). The natural candidate for an inner product, namely

$$\langle f, g \rangle := \int_{\mathcal{X}} f(x)g(x) d\mu$$

does not need to satisfy positive definiteness and it may not be a well-defined function to \mathbb{R} .

Instead, we will work in a smaller functional space $L^2(\mathcal{X})$ of all measurable functions satisfying $\int f^2(x) d\mu < +\infty$, where two functions are identified if they differ on a subset of measure zero. In this space, $\langle f, g \rangle$ defines a valid inner product.

A.2 Continuity and semicontinuity

A.2.1 Continuity

Continuity is one of the most fundamental concepts of real analysis.

Definition A.2.1. If $f : \mathcal{X} \rightarrow \mathcal{Y}$ is a mapping between metric spaces then f is **continuous** at $\mathbf{x} \in \mathcal{X}$ if for every sequence (\mathbf{x}_n) in \mathcal{X} , $\mathbf{x}_n \rightarrow \mathbf{x}$ implies $f(\mathbf{x}_n) \rightarrow f(\mathbf{x})$. We say that f is continuous if it is continuous at every \mathbf{x} .

Proposition A.2.2. Composition of continuous functions is continuous.

Proof. Follows easily from the definition. We leave the details as an exercise. \square

There are various equivalent definitions of continuity and it is important to be aware of them. The sequential definition is typically the easiest both conceptually and operationally. However, in this course we will also be using two other definitions, which we discuss next.

The standard way of defining continuous functions is, so called, (ϵ, δ) -condition.

Theorem A.2.3. A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is continuous at x if and only if

$$\forall \epsilon > 0 \exists \delta > 0 \left(d(\mathbf{x}, \mathbf{y}) < \delta \Rightarrow d(f(\mathbf{x}), f(\mathbf{y})) < \epsilon \right).$$

Another important reformulation of continuity builds on the concept of preimage. Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be given. The **preimage** of a set $V \subset \mathcal{Y}$ is

$$f^{-1}(V) := \{\mathbf{x} \in \mathcal{X} : f(\mathbf{x}) \in V\}. \quad (\text{A.1})$$

For example, if $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is defined by $f(x, y) = x^2 + y^2 + 2$ then the preimage of the interval $[3, 6]$ is the annulus in the plane with inner radius 1 and outer radius 2.

Theorem A.2.4. A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is continuous if and only if the preimage $f^{-1}(V)$ of any open set $V \subset \mathcal{Y}$ is open in \mathcal{X} .

A.2.2 Semicontinuity and Fatou's lemma

A closely related notion to continuity of a function is that of semicontinuity. Recall that $\overline{\mathbb{R}} = [-\infty, +\infty]$ is the set of extended real numbers with the usual definition of algebraic operations that incorporate the two extra "numbers" $-\infty, +\infty$. Let \mathcal{X} be any metric space.

Definition A.2.5. A function $f : \mathcal{X} \rightarrow \overline{\mathbb{R}}$ is **lower semicontinuous** at p if, for $\epsilon > 0$, p is an interior point of $\{x : f(x) > f(p) - \epsilon\}$. Similarly, f is **upper semicontinuous** at $p \in \mathcal{X}$ if, for every real $\epsilon > 0$, p is an interior point of the set $\{x : f(x) < f(p) + \epsilon\}$.

Recall that for a real sequence $(s_n)_{n \in \mathbb{N}}$ the set E of subsequential limits is the set of $s \in \mathbb{R}$ such that some subsequence of (s_n) converges to s . We define $\liminf_{n \rightarrow \infty} s_n = \inf E$ and $\limsup_{n \rightarrow \infty} s_n = \sup E$. This definition can be extended to functional limits. Consider all sequences $x_n \rightarrow p$ and the corresponding converging subsequences $f(x_{n_k})$. Then $\liminf_{x \rightarrow p} f(x)$ is the inferior of the set of all such subsequential limits. The upper functional limit $\limsup_{x \rightarrow p} f(x)$ is defined analogously.

There are important alternative ways to formulate semi-continuity.

Exercise A.2.6. Show that $f : \mathcal{X} \rightarrow \overline{\mathbb{R}}$ is lower semicontinuous at $p \in \mathcal{X}$ if and only if $\liminf_{x \rightarrow p} f(x) \geq f(p)$. Similarly, $f : \mathcal{X} \rightarrow \overline{\mathbb{R}}$ is upper semicontinuous at $p \in \mathcal{X}$ if and only if $\limsup_{x \rightarrow p} f(x) \leq f(p)$.

This exercise immediately gives us the following result.

Theorem A.2.7. A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is continuous at p if and only if it is both upper and lower semicontinuous at p .

We say that a function is lower (upper) semicontinuous if it is lower (upper) semicontinuous at every point $p \in \mathcal{X}$.

Exercise A.2.8. Show that $f : \mathcal{X} \rightarrow \overline{\mathbb{R}}$ is lower semicontinuous if and only if $\{x : f(x) > y\}$ is open for every $y \in \mathbb{R}$, or equivalently, when $\{x : f(x) \leq y\}$ is closed.

Lower semi-continuity¹ is an essential concept in convex analysis. In probability it appears, for example, in the Fatou's lemma, which we will use in later chapters.

¹ For more on semicontinuity see <https://pzwiernik.github.io/docs/RealAnalysisNotes.pdf>.

Theorem A.2.9 (Fatou's Lemma). Given a measure space $(\mathcal{X}, \mathcal{F}, \mu)$ and a set $U \in \mathcal{F}$, let $\{f_n\}$ be a sequence of measurable nonnegative functions $f_n : U \rightarrow [0, +\infty]$. Define the function $f : U \rightarrow [0, \infty]$ by setting

$$f(x) = \liminf_{n \rightarrow \infty} f_n(x) \quad \text{for all } x \in U.$$

Then f is measurable, and also

$$\int_U f d\mu \leq \liminf_{n \rightarrow \infty} \int_U f_n(x) d\mu,$$

where the integrals may be infinite.

In particular, this last statement is saying that $Z(f) = \int_{\mathcal{X}} f d\mu$ is lower semicontinuous. Namely, for every measurable f and a sequence f_n that converges pointwise to f ($\lim_n f_n(x) = f(x)$ for all $x \in \mathcal{X}$),

$$Z(f) \leq \liminf_n Z(f_n)$$

(c.f Exercise A.2.6).

A.3 Differentiation

Given a function $f : V \rightarrow \mathbb{R}$, from a vector space V , we may be interested in a local behaviour of f around some point $\mathbf{a} \in V$. The fundamental tool is given by the directional derivative. If \mathbf{u} is a vector in V , then the directional derivative of f at $\mathbf{a} \in V$ in the direction $\mathbf{u} \in V$ is

$$D_{\mathbf{u}}f(\mathbf{a}) := \lim_{t \rightarrow 0} \frac{f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a})}{t}. \quad (\text{A.2})$$

If the directional derivative $D_{\mathbf{u}}f(\mathbf{a})$ exists, it depends only on the behaviour of f in a small open neighbourhood $U \subset V$ of \mathbf{a} and so the function does not need to be defined over the whole V . The sign of the derivative has an important interpretation: if $D_{\mathbf{u}}f(\mathbf{a}) > 0$ then the value of the function increases as we move infinitesimally from \mathbf{a} in the direction \mathbf{u} .

Remark A.3.1. *In infinite dimensional spaces the directional derivative is typically called the Gateaux derivative.*

If $f : V \rightarrow \mathbb{R}$ is differentiable at \mathbf{a} the directional derivatives exist and then

$$f(\mathbf{a} + t\mathbf{u}) = f(\mathbf{a}) + t D_{\mathbf{u}}f(\mathbf{a}) + r(t\mathbf{u}),$$

where the remainder r satisfies $\lim_{t \rightarrow 0} \frac{r(t\mathbf{u})}{t} = 0$ (in other words $r(t\mathbf{u}) = o(t)$). This gives a simple ‘‘algebraic’’ way of computing the directional derivatives. Before we give some examples, we note that with a choice of an inner product on V , we get

$$D_{\mathbf{u}}f(\mathbf{a}) = \langle \nabla f(\mathbf{a}), \mathbf{u} \rangle, \quad (\text{A.3})$$

where $\nabla f(\mathbf{a}) \in V$ denotes the gradient of f at \mathbf{a} .

Exercise A.3.2. *Use the Cauchy-Schwarz inequality and (A.3) to show that $\nabla f(\mathbf{a})$ is the direction of the steepest increase of f locally around \mathbf{a} .*

Example A.3.3. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) = \mathbf{x}^\top A \mathbf{x}$ for $A \in \mathbb{S}^d$. Then*

$$f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x}) = t(\mathbf{u}^\top A \mathbf{x} + \mathbf{x}^\top A \mathbf{u}) + o(t),$$

which gives that $D_{\mathbf{u}}f(\mathbf{x}) = \mathbf{u}^\top A \mathbf{x} + \mathbf{x}^\top A \mathbf{u} = \langle 2A\mathbf{x}, \mathbf{u} \rangle$ and so $\nabla f(\mathbf{x}) = 2A\mathbf{x}$.

Example A.3.4. *Let $f : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ be defined by $f(A) = \text{tr}(A^2)$. Then*

$$f(A + tU) - f(A) = t(\text{tr}(UA) + \text{tr}(AU)) + o(t),$$

which gives that $D_U f(A) = 2\text{tr}(UA) = \langle 2A^\top, U \rangle$ and so $\nabla f(A) = 2A^\top$.

A slightly more involved matrix example is given as an exercise.

Exercise A.3.5. Consider the function $f(\Sigma) = \log \det(\Sigma)$ defined for all $\Sigma \in \mathbf{S}^m$ that are positive definite. Find the gradient of this function (it is an element of \mathbf{S}^m). Hint: Computing $D_U f(\Sigma)$ consider the eigenvalues of $\Sigma^{-1/2} U \Sigma^{-1/2}$.

We also have an infinite-dimensional example.

Example A.3.6. Let $f : L^2(\mathbb{R}) \rightarrow \mathbb{R}$ be defined by $f(\varphi) = \int \varphi^2(x) dx = \|\varphi\|^2$ then, for every $u \in L^2(\mathbb{R})$,

$$f(\varphi + tu) - f(\varphi) = 2t \int \varphi(x)u(x) dx + o(t),$$

which gives that $D_U f(\varphi) = 2 \int \varphi(x)u(x) dx = 2\langle \varphi, u \rangle$ and so $\nabla f(\varphi) = 2\varphi$.

The infinite dimensional case is extremely important in semiparametric and nonparametric statistics but it also appears in the theory of optimal statistical procedures. This may require a bit more careful treatment but we will introduce relevant concepts of functional analysis along the way. The message we tried to convey above is that very often it is useful to think about this infinite dimensional case as a special case of the standard analysis on \mathbb{R}^d .

To conclude this section we generalize (A.3) by defining the derivative as a linear map.

Definition A.3.7. Suppose U is open in \mathbb{R}^n , $f : U \rightarrow \mathbb{R}^m$. The function f is differentiable at $\mathbf{a} \in U$ with derivative f'_a if $f'_a : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear function and

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + f'_a(\mathbf{h}) + \mathbf{r}(\mathbf{h}),$$

where the remainder $\mathbf{r}(\mathbf{h})$ satisfies

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}_n} \frac{\mathbf{r}(\mathbf{h})}{\|\mathbf{h}\|} = \mathbf{0}_m.$$

If f is differentiable at every $\mathbf{a} \in U$, we say that f is differentiable in U .

Exercise A.3.8. Check if the above examples are differentiable.

Proposition A.3.9. Let $f : U \rightarrow \mathbb{R}^m$, where U is an open subset in \mathbb{R}^n , be differentiable at $\mathbf{a} \in U$. If $\mathbf{u} \in \mathbb{R}^n$ then

$$D_U f(\mathbf{a}) = f'_a(\mathbf{u}).$$

Proof. Since f'_a exists

$$f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a}) = f'_a(t\mathbf{u}) + \mathbf{r}(t\mathbf{u})$$

with $\frac{\mathbf{r}(t\mathbf{u})}{t} \rightarrow \mathbf{0}$ as $t \rightarrow 0$. Dividing by t and taking the limit, we get that

$$f'_a(\mathbf{u}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a})}{t} = D_U f(\mathbf{a}).$$

□

Exercise A.3.10. Given the same set-up as in Proposition A.3.9, show that $f'_a(\mathbf{u}) = \mathbf{J}f(\mathbf{a}) \cdot \mathbf{u}$, where $\mathbf{J}f(\mathbf{a})$ is the Jacobian of f at \mathbf{a} , that is $\mathbf{J}f(\mathbf{a}) \in \mathbb{R}^{m \times n}$ and its (i, j) -th entry is $\frac{\partial}{\partial x_j} f_i$.

B

Convex Analysis

B.1 Convexity and hyperplane separation

Let V be a vector space. Recall that a set $C \subseteq V$ is *convex* if for any two points $\mathbf{x}, \mathbf{y} \in C$ and for all $\lambda \in (0, 1)$

$$z_\lambda := (1 - \lambda)\mathbf{x} + \lambda\mathbf{y} \in C.$$

The point z_λ can be rewritten as $z_\lambda = \mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})$ so, as λ varies from 0 to 1, z_λ moves from \mathbf{x} to \mathbf{y} along the segment joining \mathbf{x} and \mathbf{y} . This gives a geometric interpretation of convex sets: for any two points in the set, the segment between these two points is contained in the set.

Given a non-empty subset $C \subset \mathbb{R}^k$ we define the **distance to C** function $d_C : \mathbb{R}^k \rightarrow \mathbb{R}$ by

$$d_C(\mathbf{x}) := \inf_{\mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|.$$

This function is well defined because for every $\mathbf{x} \in \mathbb{R}^k$ the set $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{y} \in C\} \subset \mathbb{R}$ is bounded from below (by zero) and so its infimum is well-defined. The following result is fundamental for many applications of convex analysis.

Theorem B.1.1 (Minimum distance to a set). (1) Let $E, F \subset \mathbb{R}^k$. Then

$d_F : E \rightarrow \mathbb{R}$ is Lipschitz continuous function with constant 1 (in particular, it is continuous).

(2) If F is closed, then

$$\forall \mathbf{x} \in E \quad \exists \mathbf{y} \in F \text{ such that } d_F(\mathbf{x}) = \|\mathbf{x} - \mathbf{y}\|.$$

(3) If F is also convex, for every \mathbf{x} such \mathbf{y} is unique in F .

Proof. (1) Let $\mathbf{x}_1, \mathbf{x}_2 \in E$. By the triangle inequality, $\|\mathbf{x}_1 - \mathbf{y}\| \leq \|\mathbf{x}_1 - \mathbf{x}_2\| + \|\mathbf{x}_2 - \mathbf{y}\|$. If $\mathbf{y} \in F$ then $d_F(\mathbf{x}_1) \leq \|\mathbf{x}_1 - \mathbf{y}\|$ and so

$$d_F(\mathbf{x}_1) \leq \|\mathbf{x}_1 - \mathbf{x}_2\| + \|\mathbf{x}_2 - \mathbf{y}\|.$$

Since the inequality holds for every $\mathbf{y} \in F$. Take infimum over all $\mathbf{y} \in F$ to get that $d_F(\mathbf{x}_1) \leq \|\mathbf{x}_1 - \mathbf{x}_2\| + d_F(\mathbf{x}_2)$. In the same way, starting with $\|\mathbf{x}_2 - \mathbf{y}\| \leq \|\mathbf{x}_1 - \mathbf{x}_2\| + \|\mathbf{x}_1 - \mathbf{y}\|$, conclude that $d_F(\mathbf{x}_2) \leq \|\mathbf{x}_1 - \mathbf{x}_2\| + d_F(\mathbf{x}_1)$. It follows that $|d_F(\mathbf{x}_1) - d_F(\mathbf{x}_2)| \leq \|\mathbf{x}_1 - \mathbf{x}_2\|$, which implies Lipschitz continuity of d_F .

(2) Let $\mathbf{x} \in E$ and fix $\mathbf{y}_0 \in F$. We have $d_F(\mathbf{x}) \leq \|\mathbf{x} - \mathbf{y}_0\| = r$. Define $\tilde{F} = F \cap \overline{N_r(\mathbf{x})}$, where $\overline{N_r(\mathbf{x})}$ is the closed ball of radius r around \mathbf{x} . Since F is closed and $\overline{N_r(\mathbf{x})}$ is compact, \tilde{F} is also compact. Since $\|\mathbf{x} - \mathbf{y}\|$ is a continuous function of \mathbf{y} , there exists $\mathbf{y}_1 \in \tilde{F}$ such that $\inf_{\mathbf{y} \in F} \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - \mathbf{y}_1\|$.

(3) Let \mathbf{y}_1 and \mathbf{y}_2 in F be such that $\|\mathbf{x} - \mathbf{y}_1\| = \|\mathbf{x} - \mathbf{y}_2\| = d_F(\mathbf{x})$. Define $\mathbf{p} = \mathbf{y}_2 - \mathbf{y}_1$ and $h : [0, 1] \rightarrow \mathbb{R}$ by $h(\lambda) = \|\mathbf{x} - \mathbf{y}_1 - \lambda\mathbf{p}\|^2$. We have $h(0) = h(1)$ and also, because F is convex, $\mathbf{y}_1 + \lambda\mathbf{p} \in F$ for $\lambda \in [0, 1]$ and so h is minimized at $\lambda = 0$ and $\lambda = 1$. Since $h(\lambda)$ is a quadratic function with nonnegative coefficient $\|\mathbf{y}_1 - \mathbf{y}_2\|^2$ of λ^2 , this is only possible if $\mathbf{y}_1 = \mathbf{y}_2$. \square

Proposition B.1.2. Let $E, F \subset \mathbb{R}^k$ with F closed and convex. Let $g : E \rightarrow F$ be given by $g(\mathbf{x}) = \arg \inf_{\mathbf{y} \in F} \|\mathbf{x} - \mathbf{y}\|$. Then g is a well-defined and

$$\|g(\mathbf{x}_2) - g(\mathbf{x}_1)\| \leq \|\mathbf{x}_2 - \mathbf{x}_1\| \quad \text{for all } \mathbf{x}_1, \mathbf{x}_2 \in E. \quad (\text{B.1})$$

In particular, g is a continuous function.

Proof. Because F is closed and convex, Theorem B.1.1 assures that g is a well-defined function, that is, for each $\mathbf{x} \in E$ there is a unique $\mathbf{y} \in F$ such that $g(\mathbf{x}) = \mathbf{y}$. To show (B.1), take $\mathbf{p} = g(\mathbf{x}_2) - g(\mathbf{x}_1)$;

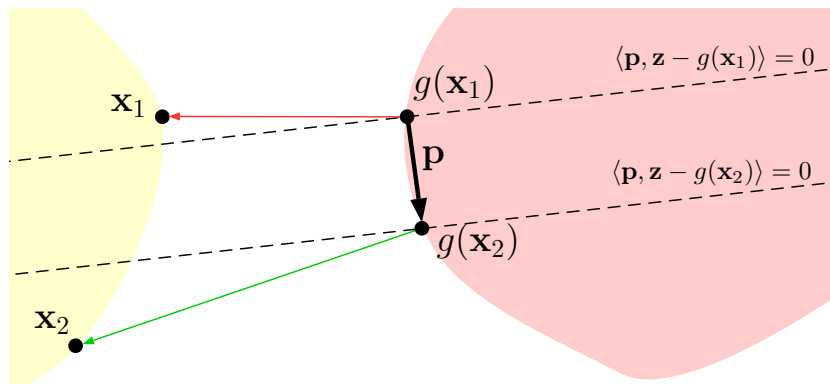


Figure B.1: Illustration of the proof of Proposition B.1.2

c.f. Figure B.1. The set of points $g(\mathbf{x}_1) + t\mathbf{p}$ for $t \in [0, 1]$ lies in F , by convexity, and so

$$h(t) = \|\mathbf{x}_1 - g(\mathbf{x}_1) - t\mathbf{p}\|^2$$

has a minimum at $t = 0$. This is a quadratic function in t with a strictly positive coefficient of t^2 . The only way for such a function to

have a minimum at $t = 0$ is that its derivative at $t = 0$ is nonnegative, or, in other words, the coefficient of t is nonnegative. This coefficient is $-2\langle \mathbf{p}, \mathbf{x}_1 - g(\mathbf{x}_1) \rangle$, which implies that $\langle \mathbf{p}, \mathbf{x}_1 - g(\mathbf{x}_1) \rangle \leq 0$. In a similar way, we show that $\langle \mathbf{p}, \mathbf{x}_2 - g(\mathbf{x}_2) \rangle \geq 0$. But these two inequalities imply that

$$\langle \mathbf{p}, \mathbf{x}_2 - \mathbf{x}_1 \rangle \geq \langle \mathbf{p}, g(\mathbf{x}_2) - g(\mathbf{x}_1) \rangle = \|\mathbf{p}\|^2.$$

The Cauchy-Schwarz inequality gives that $\langle \mathbf{p}, \mathbf{x}_2 - \mathbf{x}_1 \rangle \leq \|\mathbf{p}\| \|\mathbf{x}_2 - \mathbf{x}_1\|$. This implies that $\|\mathbf{p}\| \leq \|\mathbf{x}_2 - \mathbf{x}_1\|$, which is precisely (B.1). \square

We say that two sets $A, B \subset \mathbb{R}^k$ are **separated by a hyperplane** if there exists $\mathbf{p} \neq \mathbf{0}_k$ and $c \in \mathbb{R}$ such that for every $\mathbf{x} \in A, \mathbf{y} \in B$ we have $\langle \mathbf{p}, \mathbf{x} \rangle \leq c \leq \langle \mathbf{p}, \mathbf{y} \rangle$. The separation is strict if the inequalities can be made strict. The following theorem is one of the most important results of convex geometry.

Theorem B.1.3 (Hyperplane Separation Theorem). *Let C and K be disjoint non-empty convex sets in \mathbb{R}^k . Then they are separated by a hyperplane. Let C be closed and K compact. Then C and K are strictly separated by a hyperplane.*

Note that compactness of K in the second part of Theorem B.1.3 is necessary. For example $A = \{\mathbf{x} \in \mathbb{R}^2 : x_1 \leq 0\}$, $B = \{\mathbf{x} \in \mathbb{R}^2 : x_1 > 0, x_2 \geq 1/x_1\}$ are closed but not strictly separated.

Proof. We prove the second part of the theorem leaving the first part as an exercise. By Theorem B.1.1, $d_C(\mathbf{x})$ is a continuous real-valued function on K and so it achieves its minimum. Call it $\mathbf{x}_0 \in K$. By Theorem B.1.1(c) there is exactly one $\mathbf{y}_0 \in C$ such that $d_C(\mathbf{x}_0) = \|\mathbf{x}_0 - \mathbf{y}_0\|$. Set $\mathbf{p} = \mathbf{x}_0 - \mathbf{y}_0$. Then $\mathbf{p} \neq \mathbf{0}_k$ and $0 < \|\mathbf{p}\|^2 = \langle \mathbf{p}, \mathbf{x}_0 - \mathbf{y}_0 \rangle$ so

$$\langle \mathbf{p}, \mathbf{x}_0 \rangle > \langle \mathbf{p}, \mathbf{y}_0 \rangle$$

so it suffices to show that $\langle \mathbf{p}, \mathbf{x} \rangle \geq \langle \mathbf{p}, \mathbf{x}_0 \rangle$ for every $\mathbf{x} \in K$ and $\langle \mathbf{p}, \mathbf{y} \rangle \leq \langle \mathbf{p}, \mathbf{y}_0 \rangle$ for all $\mathbf{y} \in C$. We show the second, the first is similar and follows from convexity. Let $\mathbf{y} \in C$ and set $\mathbf{y}_\lambda = (1 - \lambda)\mathbf{y}_0 + \lambda\mathbf{y}$. Then

$$\mathbf{x}_0 - \mathbf{y}_\lambda = \mathbf{x}_0 - \mathbf{y}_0 - \lambda(\mathbf{y} - \mathbf{y}_0) = \mathbf{p} - \lambda(\mathbf{y} - \mathbf{y}_0)$$

and so

$$\begin{aligned} \|\mathbf{x}_0 - \mathbf{y}_\lambda\|^2 &= \|\mathbf{p} - \lambda(\mathbf{y} - \mathbf{y}_0)\|^2 = \\ &= \lambda^2\|\mathbf{y} - \mathbf{y}_0\|^2 - 2\lambda\langle \mathbf{p}, \mathbf{y} - \mathbf{y}_0 \rangle + \|\mathbf{p}\|^2. \end{aligned}$$

This is a quadratic function of λ that achieves its minimum at $\lambda = 0$ (by construction!). This implies that the derivative of this function

at zero must be nonnegative. This derivative is equal to the coefficient of λ which is $2\langle \mathbf{p}, \mathbf{y}_0 - \mathbf{y} \rangle$. This implies that $\langle \mathbf{p}, \mathbf{y} \rangle \leq \langle \mathbf{p}, \mathbf{y}_0 \rangle$ as claimed. \square

Exercise B.1.4. Show that each closed convex set is an intersection of closed half-spaces.

B.2 Convex functions and optimization

A function f defined on a convex set $C \subseteq V$ with values in $\overline{\mathbb{R}}$ is convex if

$$f((1-\lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1-\lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) \quad \text{for all } \mathbf{x} \neq \mathbf{y} \in C, \lambda \in (0,1).$$

Moreover, f is strictly convex if the inequality is always strict. A function f is (strictly) concave if $-f$ is (strictly) convex.

Exercise B.2.1. Let \mathbb{S}_+^m be the set of all symmetric $m \times m$ that are positive definite. Show that the following functions are convex:

- (a) $f : \mathbb{R}^d \rightarrow \mathbb{R}$ defined by $f(\mathbf{x}) = \|\mathbf{x}\|^2$.
- (b) $f : \mathbb{S}_+^m \rightarrow \mathbb{R}$ defined by $f(\Sigma) = -\log \det(\Sigma)$.
- (c) $L : L^2(\mathbb{R}) \rightarrow \mathbb{R}$ defined by $L(f) = \int_{-\infty}^{\infty} f^2(x) dx$.

Are they strictly convex?

We now discuss the most important features of convex function. Note that we never assume that V is a finite-dimensional space.

Proposition B.2.2. If $f : C \rightarrow \overline{\mathbb{R}}$ is convex and $\mathcal{L} \subset V$ is a linear subspace then f restricted to $C \cap \mathcal{L}$ is also convex.

Proof. This follows directly from the definition. \square

Theorem B.2.3 (Jensen's inequality). Suppose $f : C \rightarrow \mathbb{R}$ is a convex function and let X be a random variable with $\mathbb{P}(X \in C) = 1$ and $\mathbb{E}X < \infty$. Then $f(\mathbb{E}(X)) \leq \mathbb{E}(f(X))$. If f is strictly convex then the inequality is strict unless X is constant almost surely.

Proposition B.2.4. If \mathbf{a} is a local optimum of a convex function f , then f is a global optimum.

Proof. We argue by contradiction. If \mathbf{x}, \mathbf{y} are two local optima with $f(\mathbf{x}) < f(\mathbf{y})$ (\mathbf{y} is local but not global) then, for every $\lambda \in (0,1)$, $f(z_\lambda) < f(\mathbf{y})$. This contradicts local optimality of \mathbf{y} . \square

Proposition B.2.5. Let $g : C \times \mathcal{A} \rightarrow \mathbb{R}$ be such that $g(x, \alpha)$ is convex in $x \in C$ for every fixed α . Then the function

$$f(x) := \sup_{\alpha \in \mathcal{A}} g(x, \alpha)$$

(defined as the pointwise supremum) is a convex function.

Proof. Exercise. □

Example B.2.6 (Fenchel conjugate). Let $f(x)$ be a convex function on C . Then the function

$$g(x, y) = \langle x, y \rangle - f(x)$$

is concave in x and linear (and so also convex) in y . Using Proposition B.2.5 we get

$$f^*(y) = \sup_{x \in C} \{\langle x, y \rangle - f(x)\}$$

is a convex function. The function f^* is called the Fenchel conjugate of f . Note that the Fenchel inequality follows easily:

$$f(x) + f^*(y) \geq \langle x, y \rangle \quad \text{for all } x, y.$$

For a simple example, take $f(x) = x^2$.

Proposition B.2.4 suggests that optimizing convex functions is much easier than optimizing general functions. However, the study of local behaviour of convex functions is easy also for a different reason. Define a one-sided directional derivative

$$D_u^+ f(a) = \lim_{t \rightarrow 0^+} \frac{f(a + tu) - f(a)}{t}. \quad (\text{B.2})$$

If the directional derivative $D_u f(a)$ in (A.2) exists then the one-sided derivative exists and they are equal. But there are important examples when the directional derivative does not exist but the one-sided derivative does.

Exercise B.2.7. Show that at the origin none of the directional derivatives of $f(\mathbf{x}) = \|\mathbf{x}\|_1$ exists but all the one-sided derivatives do.

The importance of one-sided directional derivatives comes from the fact that if $D_u^+ f(a)$ is positive then the function is increasing when we move infinitesimally from a in the direction u . This can be used to easily provide necessary conditions for a local optimum even in the constrained setting.

Theorem B.2.8. Let $f : C \rightarrow \mathbb{R} \cup \{\infty\}$ be convex, and let $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$. Then for every $\mathbf{z}_\lambda = (1 - \lambda)\mathbf{x} + \lambda\mathbf{y}$, $\lambda \in (0, 1)$, the one-sided derivative $D_{\mathbf{x}-\mathbf{y}}^+ f(\mathbf{z}_\lambda)$ exists and is an increasing function of λ .

We claim that this result is true also over infinite-dimensional spaces. Our proof of Theorem B.2.8 starts with the following lemma.

Lemma B.2.9. Let $\varphi : [0, 1] \rightarrow \mathbb{R}$ be convex. Consider the “chords” $0 \leq p < q < 1$, and $0 < s < t \leq 1$, with $p \leq s$ and $q \leq t$. Then

$$\frac{\varphi(q) - \varphi(p)}{q - p} \leq \frac{\varphi(t) - \varphi(s)}{t - s}.$$

Proof. Since $p < q \leq t$, then

$$q = \frac{t-q}{t-p}p + \frac{q-p}{t-p}t$$

is a convex combination of p and t , so the convexity of φ gives

$$\varphi(q) \leq \frac{t-q}{t-p}\varphi(p) + \frac{q-p}{t-p}\varphi(t).$$

This is equivalent to

$$\frac{\varphi(q) - \varphi(p)}{q-p} \leq \frac{\varphi(t) - \varphi(p)}{t-p}. \quad (\text{B.3})$$

Now, apply this inequality to the function $\psi(\lambda) = \varphi(1-\lambda)$ and the points $1-t < 1-s \leq 1-p$. The result may be rewritten as

$$\frac{\varphi(t) - \varphi(p)}{t-p} \leq \frac{\varphi(t) - \varphi(s)}{t-s}.$$

Together with (B.3), this is the required inequality. \square

Proof of Theorem B.2.8. Let $\varphi(\lambda) = f(\mathbf{z}_\lambda)$ and apply Lemma B.2.9 for the points $p = s = \lambda < q \leq t$ to obtain

$$\frac{\varphi(q) - \varphi(\lambda)}{q-\lambda} \leq \frac{\varphi(t) - \varphi(\lambda)}{t-\lambda}.$$

It follows that $\psi(q) = (\varphi(q) - \varphi(\lambda))/(q - \lambda)$ is an increasing function of q , $q > \lambda$. Note that

$$\psi(q) = t^{-1}\{f(\mathbf{z}_\lambda + t(\mathbf{x} - \mathbf{y})) - f(\mathbf{z}_\lambda)\}, \quad t = q - \lambda.$$

Applying Lemma B.2.9 to the points $p < s = q = \lambda < t$ shows that

$$\psi(p) = \frac{\varphi(\lambda) - \varphi(p)}{\lambda-p} \leq \frac{\varphi(t) - \varphi(\lambda)}{t-\lambda} = \psi(t).$$

Hence, $\psi(t)$ is increasing and bounded below to $t > \lambda$, and thus,

$$\lim_{t \rightarrow \lambda^+} \psi(t)$$

exists. Equivalently the one-sided derivative $D_{\mathbf{x}-\mathbf{y}}^+(\mathbf{z}_\lambda)$ exists for $\lambda \in (0, 1)$. By taking limits in the inequality in Lemma B.2.9, namely $q \rightarrow p^+$, $t \rightarrow s^+$, we see that $\varphi'(p) \leq \varphi'(s)$ (one-sided derivatives) for almost all p and s with $p \leq s$. Thus $D_{\mathbf{x}-\mathbf{y}}^+(\mathbf{z}_\lambda)$ is an increasing function of λ . \square

Exercise B.2.10. Show that any convex function defined and finite on a convex set C must be continuous on its interior. Although this result is general, for simplicity, you can focus on its one-dimensional version.

We now informally state an important result, which gives a fundamental understanding of a large family of optimization problems. Often we optimize a function over a convex set C given by bunch of linear constraints together with inequality constraints $g_i(x) \leq 0$ for $i \in \mathcal{I}$ with g_i convex. If the functions g_i are all continuously differentiable and all one-sided derivatives in (B.2) exist then \mathbf{a} is a local minimum of f if and only if $D_{\mathbf{u}}^+ f(\mathbf{a}) \geq 0$ for all \mathbf{u} such that $D_{\mathbf{u}} g_i(\mathbf{a}) \leq 0$ for all i such that $g_i(\mathbf{a}) = 0$ (active constraints).

C

Probability

C.1 Continuity of probability

In what follows we fix a probability space $(\Omega, \mathcal{B}, \mathbb{P})$. If A_n is a sequence of events then we say that (A_n) increases to A if $A_n \subseteq A_{n+1}$ for all $n \in \mathbb{N}$ and $A = \bigcup_{n \geq 1} A_n$.

Proposition C.1.1. *If (A_n) is a sequence of events increasing to A , then*

$$\lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}(A).$$

Proof. Clearly $A_n = A_n \cap A_{n+1}$ for all $n \in \mathbb{N}$. Let $A_0 = \emptyset$, $C_1 = A_1$, and define $C_{n+1} = A_{n+1} \setminus A_n$. Notice that C_1, C_2, \dots are disjoint with

$$\bigcup_{j=1}^n C_j = A_n \quad \text{and} \quad \bigcup_{j=1}^{\infty} C_j = \bigcup_{j=1}^{\infty} A_j = A.$$

Using the fact that \mathbb{P} is countably additive, we conclude

$$\mathbb{P}(A) = \mathbb{P}\left(\bigcup_{j=1}^{\infty} C_j\right) = \sum_{j=1}^{\infty} \mathbb{P}(C_j) = \lim_{n \rightarrow \infty} \sum_{j=1}^n \mathbb{P}(C_j) = \lim_{n \rightarrow \infty} \mathbb{P}\left(\bigcup_{j=1}^n C_j\right) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n).$$

□

This immediately gives the following result.

Proposition C.1.2. *Both the CDF $F(t) = \mathbb{P}(X \leq t)$ of X and its survival function $G(t) = \mathbb{P}(X > t)$ are right-continuous.*

Proof. Let (t_n) be any monotone sequence such that $t_n > t$ and $t_n \rightarrow t$. Then the events $A_n = \{X > t_n\}$ satisfy $A_n \subseteq A_{n+1}$ for all $n \in \mathbb{N}$. Since $\bigcup_{n \in \mathbb{N}} A_n = \{X > t\}$ we have that (A_n) increases to $\{X > t\}$. By Proposition C.1.1,

$$\lim_{n \rightarrow \infty} \mathbb{P}(X > t_n) = \mathbb{P}(X > t).$$

Since t_n is otherwise arbitrary, we conclude $\lim_{x \rightarrow t^+} G(x) = G(t)$ and so the survival function is right-continuous. The claim for the CDF follows immediately as

$$\lim_{x \rightarrow t^+} F(x) = \lim_{x \rightarrow t^+} (1 - G(x)) = 1 - \lim_{x \rightarrow t^+} G(x) = 1 - G(t) = F(t).$$

□

Remark C.1.3. Using the same approach we can show that the functions $\mathbb{P}(X < t)$ and $\mathbb{P}(X \geq t)$ are left-continuous. We leave it as an exercise.

C.2 Martingales

Let X_1, \dots, X_n be independent random variables with values in \mathcal{X} . Let $\{\mathcal{F}_k\}_{k=1}^\infty$ be a sequence of σ -fields, $\mathcal{F}_k = \sigma(X_1, \dots, X_k)$. Let $\{Y_k\}_{k=1}^\infty$ be a sequence of variables such that Y_k is \mathcal{F}_k -measurable (we say that $\{Y_k\}_{k=1}^\infty$ is adapted to the filtration $\mathcal{F}_k = \sigma(X_1, \dots, X_k)$).

Definition C.2.1. Given a sequence $\{Y_k\}_{k=1}^\infty$ adapted to the filtration $\mathcal{F}_k = \sigma(X_1, \dots, X_k)$, the pair $\{(Y_k, \mathcal{F}_k)\}_{k=1}^\infty$ is a **martingale** if, for all $k \geq 1$,

$$\mathbb{E}(|Y_k|) < \infty \quad \text{and} \quad \mathbb{E}[Y_{k+1} | \mathcal{F}_k] = Y_k.$$

Example C.2.2 (Simple random walk). A particle jumps either one step to the right or one step to the left with the corresponding probabilities p and $q = 1 - p$. Assume that the subsequent moves are independent of each other. Define $S_n = X_1 + \dots + X_n$. It is clear that $\mathbb{E}|S_n| \leq n$ and

$$\mathbb{E}[S_{n+1} | X_1, \dots, X_n] = S_n + (p - q),$$

and so $Y_n = S_n - n(p - q)$ defined a martingale with respect to X .

Example C.2.3 (Likelihood ratio). Let p_1, p_0 be two mutually absolutely continuous densities, and let X_1, X_2, \dots be an i.i.d. sequence from p_0 . For each $k \in \mathbb{N}$ let $Y_k = \prod_{i=1}^k \frac{p_1(X_i)}{p_0(X_i)}$ be the likelihood ratio based on the first k samples (c.f (3.15)). Then the sequence is a martingale with respect to $\{X_k\}_{k=1}^\infty$. Indeed,

$$\mathbb{E}[Y_{k+1} | \mathcal{F}_k] = \mathbb{E} \left[\frac{p_1(X_{k+1})}{p_0(X_{k+1})} \prod_{i=1}^k \frac{p_1(X_i)}{p_0(X_i)} \right] = Y_k,$$

using the fact that $\mathbb{E} \left[\frac{p_1(X_{k+1})}{p_0(X_{k+1})} \right] = 1$.

There are many cases of interest in which the martingale condition $\mathbb{E}(Y_{k+1} | \mathcal{F}_k) = Y_k$ does not hold, being replaced instead by an inequality: $\mathbb{E}(Y_{k+1} | \mathcal{F}_k) \geq Y_k$ for all k , or by $\mathbb{E}(Y_{k+1} | \mathcal{F}_k) \leq Y_k$ for all k . Sequences satisfying such inequalities have many of the properties of martingales. Recall $x^+ = \max\{0, x\}$ and $x^- = -\min\{0, x\}$.

Definition C.2.4. Given a sequence $\{Y_k\}_{k=1}^\infty$ adapted to the filtration $\mathcal{F}_k = \sigma(X_1, \dots, X_k)$, the pair $\{(Y_k, \mathcal{F}_k)\}_{k=1}^\infty$ is a **submartingale** if, for all $k \geq 1$,

$$\mathbb{E}(Y_k^+) < \infty \quad \text{and} \quad \mathbb{E}[Y_{k+1} | \mathcal{F}_k] \geq Y_k,$$

or a **supermartingale** if, for all $k \geq 1$,

$$\mathbb{E}(Y_k^-) < \infty \quad \text{and} \quad \mathbb{E}[Y_{k+1} | \mathcal{F}_k] \leq Y_k,$$

We call the pair $\{(S_k, \mathcal{F}_k)\}_{k=1}^\infty$ predictable if S_{k+1} is \mathcal{F}_k -measurable for all k . We call a predictable process $\{(S_k, \mathcal{F}_k)\}_{k=1}^\infty$ increasing if $S_1 = 0$ and $\mathbb{P}(S_{k+1} \geq S_k) = 1$ for all k .

Theorem C.2.5 (Doob decomposition). A submartingale Y_k with finite means may be expressed in the form

$$Y_k = M_k + S_k,$$

where M_k is a martingale and S_k is an increasing predictable process. This decomposition is unique.

A closely related notion to martingales is that of a **martingale difference sequence**, which is an adapted sequence $\{\Delta_k, \mathcal{F}_k\}_{k=1}^\infty$ such that for all $k \geq 1$,

$$\mathbb{E}|\Delta_k| \leq \infty \quad \text{and} \quad \mathbb{E}(\Delta_{k+1} | \mathcal{F}_k) = 0.$$

If $\{Y_k\}$ is a martingale then $\Delta_k = Y_k - Y_{k-1}$ is a martingale difference sequence. In our case, this easily follows from the fact that $\mathbb{E}_k(\mathbb{E}_{k+1}(Z)) = \mathbb{E}_k Z$.

