# Introduction to Machine Learning

Laurent Younes

September 5, 2024

# Contents

# Preface

Machine learning addresses the issue of analyzing, reproducing and predicting various mechanisms and processes observable through experiments and data acquisition. With the impetus of large technological companies in need of leveraging information included in the gigantic datasets that they produced or obtained through user data, with the development of new data acquisition techniques in biology, physics or astronomy, with the improvement of storage capacity and high-performance computing, this field has experienced an explosive growth over the past decades, in terms of scientific production and technological impact.

While it is being recognized in some places as a scientific discipline in itself, machine learning (which has received a few almost synonymic denominations across time, including artificial intelligence, machine intelligence or statistical learning), can also be seen as an interdisciplinary field interfacing techniques from traditional domains such as computer science, applied mathematics, and statistics. From statistics, and more specially nonparametric statistics, it borrows its main formalism, asymptotic results and generalization bounds. It also builds on many classical methods that have been developed for estimation and prediction. From computer science, it involves the construction and implementation of efficient algorithms, programming design and architecture. Finally, machine learning leverages classical methods from linear algebra and functional analysis, as well as from convex and nonlinear optimization, fields within which it had also provided new problems and discoveries. It forms a significant part of the larger field commonly called "data science," which includes methods for storing, sharing and managing data, the development powerful computer architectures for increasingly demanding algorithms, and, importantly, the definition of ethical limits and processes through which data should be used in the modern world.

This book, which originates from lecture notes of a series of graduate course taught in the Department of Applied Mathematics and Statistics at Johns Hopkins University, adopts a viewpoint (or bias) mainly focused on the mathematical and statistical aspects of the subject. Its goal is to introduce the mathematical foundations and techniques that lead to the development and analysis of many of the algorithms that are used today. It is written with the hope to provide the reader with a deeper

understanding of the algorithms made available to her in multiple machine learning packages and software, and that she will be able to assess their prerequisites and limitations, and to extend them and develop new algorithms. Note that, while adopting a presentation with a strong mathematical flavor, we will still make explicit the details of many important machine learning algorithms.

Unsurprisingly, the book will be more accessible to a reader with some background in mathematics and statistics. It assumes familiarity with basic concepts in linear algebra and matrix analysis, in multivariate calculus and in probability and statistics. We tried to place a limit at the use of measure theoretic tools, that are avoided up to a few exceptions, which are be localized and be accompanied with alternative interpretations allowing for a reading at a more elementary level.

The book starts with an introductory chapter that describes notation used throughout the book and serve at a reminder of basic concepts in calculus, linear algebra and probability. It also introduces some measure theoretic terminology, and can be used as a reading guide for the sections that use these tools. This chapter is followed by two chapters offering background material on matrix analysis and optimization. The latter chapter, which is relatively long, provides necessary references to many algorithms that are used in the book, including stochastic gradient descent, proximal methods, etc.

Chapter 4, which is also introductory, illustrates the bias-variance dilemma in machine learning through the angle of density estimation and motivates chapter 5 in which basic concepts for statistical prediction are provided. Chapter 6 provides an introduction to reproducing kernel theory and Hilbert space techniques that are used in many places, before tackling, with chapters 7 to 11, the description of various algorithms for supervised statistical learning, including linear methods, support vector machines, decision trees, boosting, or neural networks.

Chapter 12, which presents sampling methods and an introduction to the theory of Markov chains, starts a series of chapters on generative models, and associated learning algorithms. Graphical models and described in chapters 13 to 15. Chapter 16 introduces variational methods for models with latent variables, with applications to graphical models in chapter 17. Generative techniques using deep learning are presented in chapter 18.

Chapters 19 to 21 focus on unsupervised learning methods, for clustering, factor analysis and manifold learning. The final chapter of the book is theory-oriented and discusses concentration inequalities and generalization bounds.

# Chapter 1

# General Notation and Background Material

## 1.1   Linear algebra

**1.** The set of all subsets of a given set $A$ is denoted $\mathcal{P}(A)$. If $A$ and $B$ are two sets, the notation $B^A$ refers to the set of all functions $f : A \to B$. In particular, $\mathbb{R}^A$ is the space of real-valued functions, and forms a vector space. When $A$ is finite, this space is finite dimensional and can be identified with $\mathbb{R}^{|A|}$, where $|A|$ denotes the cardinality (number of elements) of $A$.

The indicator function of a subset $C$ of $A$ will be denoted $\mathbf{1}_C : A \to \{0, 1\}$, with $\mathbf{1}_C(x) = 1$ if $x \in C$ and 0 otherwise. We will sometimes write $\mathbf{1}_{x \in C}$ for $\mathbf{1}_C(x)$.

**2.** Elements of the $d$-dimensional Euclidean space $\mathbb{R}^d$ will be denoted with letters such as $x, y, z$, and their coordinates will be indexed as parenthesized exponents, so that

$$x = \begin{pmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{pmatrix}$$

(we will always identify element of $\mathbb{R}^d$ with column vectors). We will not distinguish in the notation between "points" in $\mathbb{R}^d$, seen as an affine space, and "vectors" in $\mathbb{R}^d$, seen as a vector space. The vectors $\mathbf{0}_d$ and $\mathbb{1}_d$ will denote the $d$-dimensional vectors with all coordinates equal to 0 and 1, respectively. The identity matrix in $\mathbb{R}^d$ will be denoted $\mathrm{Id}_{\mathbb{R}^d}$. The canonical basis of $\mathbb{R}^d$, provided by the columns of $\mathrm{Id}_{\mathbb{R}^d}$ will be denoted $\mathfrak{e}_1, \ldots, \mathfrak{e}_d$.

**3.** The Euclidean norm of a vector $x \in \mathbb{R}^d$ is denoted $|x|$ with

$$|x| = \left( (x^{(1)})^2 + \cdots + (x^{(d)})^2 \right)^{1/2}.$$

It will sometimes be denoted $|x|_2$, identifying it as a member of the family of $\ell^p$ norms

$$|x|_p = \left( (x^{(1)})^p + \cdots + (x^{(d)})^p \right)^{1/p} \tag{1.1}$$

for $p \geq 1$. One can also define $|x|_p$ for $0 < p < 1$, using (1.1), but in this case one does not get a norm because the triangle inequality $|x + y|_p \leq |x|_p + |y|_p$ is not true in general. The family is interesting, however, because it approximates, in the limit $p \to 0$, the number of non-zero components of $x$, denoted $|x|_0$, which is a measure of sparsity. Note that we also use the notation $|A|$ to denote the cardinality (number of elements) of a set $A$, hopefully without risk of confusion.

While we use single bars ($|x|$) to represent norms of finite-dimensional vectors, we will use double bars ($\|h\|$) for infinite-dimensional objects.

**4.** The set of $m \times d$ real matrices with real entries is denoted $\mathcal{M}_{m,d}(\mathbb{R})$, or simply $\mathcal{M}_{m,d}$ ($\mathcal{M}_{d,d}$ will also be denoted $\mathcal{M}_d$). The set of invertible $d \times d$ matrices will be denote $\mathcal{GL}_d(\mathbb{R})$.

Given $m$ column vectors $x_1, \ldots, x_m \in \mathbb{R}^d$, the notation $[x_1, \ldots, x_m]$ refers to the $d$ by $m$ matrix with $j^{\text{th}}$ column equal to $x_j$, so that, for example, $\text{Id}_{\mathbb{R}^d} = [\mathfrak{e}_1, \ldots, \mathfrak{e}_d]$.

Entry $(i, j)$ in a matrix $A \in \mathcal{M}_{m,d}(\mathbb{R})$ will either be denoted $A(i, j)$ or $A_j^{(i)}$. The rows of $A$ will be denoted $A^{(1)}, \ldots, A^{(m)}$ and the columns $A_1, \ldots, A_m$.

The operator norm of a matrix $A \in \mathcal{M}_{m,d}$ is defined by

$$|A|_{\text{op}} = \max\{|Ax| : x \in \mathbb{R}^d, |x| = 1\}.$$

**5.** The space of $d \times d$ real symmetric matrices is denoted $\mathcal{S}_d$, and its subsets containing positive semi-definite (resp. positive definite) matrices is denoted $\mathcal{S}_d^+$ (resp. $\mathcal{S}_d^{++}$). If $m \leq d$, $\mathcal{O}_{m,d}$ denotes the set of $m \times d$ matrices $A$ such that $AA^T = \text{Id}_{\mathbb{R}^m}$, and one writes $\mathcal{O}_d$ for $\mathcal{O}_{d,d}$, the space of $d$-dimensional orthogonal matrices. Finally, $\mathcal{SO}_d$ is the subset $\mathcal{O}_d$ containing orthogonal matrices with determinant 1, i.e., rotation matrices.

**6.** A $k$-multilinear mapping is a function $a : (x_1, \ldots, x_k) \mapsto a(x_1, \ldots, x_k)$ defined on $(\mathbb{R}^d)^k$ with values in $\mathbb{R}^q$ which is linear in each of its variables. The mapping is symmetric if its value is unchanged after any permutation of the variables. If $k = 2$ and $q = 1$, one also says that $a$ is a bilinear form. The norm of a $k$-multilinear mapping is defined as

$$|a| = \max\{a(x_1, \ldots, x_k) : |x_j| \leq 1, j = 1, \ldots, k\}$$

so that

$$|a(x_1, \ldots, x_k)| \leq |a| \prod_{j=1}^{k} |x_j|$$

for all $x_1, \ldots, x_k \in \mathbb{R}^d$.

A symmetric bilinear form $a$ is called positive semidefinite if $a(x, x) \geq 0$ for all $x \in \mathbb{R}^d$, and positive definite if it is positive semi-definite and $a(x, x) = 0$ if and only if $x = 0$.

Symmetric bilinear forms can always be expressed in the form $a(x,y) = x^T A y$ for some symmetric matrix $A$, and $a$ is positive (semi-)definite if and only $A$ is also. Analogous statements hold for negative (semi-)definite forms and matrices. We will use the notation $A \succ 0$ (resp. $\succeq 0$) to indicate that $A$ is positive definite (resp. positive semidefinite). Note that, if $a(x,y) = x^T A y$ for $A \in \mathcal{S}_d$, then $|a| = |A|_{\mathrm{op}}$.

## 1.2 Topology

**1.** The open balls in $\mathbb{R}^d$ will be denoted

$$B(x,r) = \{y \in \mathbb{R}^d : |y - x| < r\},$$

with $x \in \mathbb{R}^d$ and $r > 0$. The closed balls are denoted $\bar{B}(x,r)$ and contain all $y$'s such that $|y - x| \le r$. A set $U \subset \mathbb{R}^d$ is open if and only if for any $x \in U$, there exists $r > 0$ such that $B(x,r) \subset U$. A set $\Gamma \subset \mathbb{R}^d$ is closed if its complement, denoted

$$\Gamma^c = \{x \in \mathbb{R}^d : x \notin \Gamma\}$$

is open. The topological interior of a set $A \subset \mathbb{R}^d$ is the largest open set included in $A$. It will be denoted either by $\mathring{A}$ or $\mathrm{int}(A)$. A point $x$ belongs to $\mathring{A}$ if and only if $B(x,r) \subset A$ for some $r > 0$.

**2.** The closure of $A$ is the smallest closed set that contains $A$ and will be denoted either $\bar{A}$ or $\mathrm{cl}(A)$. A point $x$ belongs to $\bar{A}$ if and only if $B(x,r) \cap A \ne \emptyset$ for all $r > 0$. Alternatively, $x$ belongs to $\bar{A}$ if and only if there exists a sequence $(x_k)$ that converges to $x$ with $x_k \in A$ for all $k$.

**3.** A compact set in $\mathbb{R}^d$ is a set $\Gamma$ such that any sequence of points in $\Gamma$ contains a subsequence that converges to some point in $\Gamma$. An alternate definition is that, whenever $\Gamma$ is covered by a collection of open sets, there exists a finite subcollection that still covers $\Gamma$.

One can show that compact subsets of $\mathbb{R}^d$ are exactly its bounded and closed subsets.

**4.** A metric space is a space $\mathcal{B}$ equipped with a distance, i.e., a function $\rho : \mathcal{B} \times \mathcal{B} \to [0, +\infty)$ that satisfies the following three properties.

$$\forall x, y \in \mathcal{B} : \rho(x,y) = 0 \Leftrightarrow x = y, \tag{1.2a}$$

$$\forall x, y \in \mathcal{B} : \rho(x,y) = \rho(y,x), \tag{1.2b}$$

$$\forall x, y, z \in \mathcal{B} : \rho(x,z) \le \rho(x,y) + \rho(y,z). \tag{1.2c}$$

Equation (1.2c) is called the triangle inequality. The norm of the difference between two points: $\rho(x,y) = |x - y|$, is a distance on $\mathbb{R}^d$. The definition of open and closed subsets in metric spaces is the same as above, with $\rho(x,y)$ replacing $|x - y|$, and one says that $(x_n)$ converges to $x$ if and only if $\rho(x_n, x) \to 0$.

Compact subsets are also defined in the same way, but are not necessarily characterized as bounded and closed.

## 1.3  Calculus

**1.** If $x, y \in \mathbb{R}^d$, we will denote by $[x, y]$ the closed segment delimited by $x$ and $y$, i.e., the set of all points $(1 - t)x + ty$ for $0 \le t \le 1$. One denotes by $[x, y)$, $(x, y]$ and $(x, y)$ the semi-open or open segments, with appropriate strict inequality for $t$. (Similarly to the notation for open intervals, whether $(x, y)$ denotes an open segment or a pair of points will always be clear from the context.)

**2.** The derivative of a differentiable function $f : t \mapsto f(t)$ from an interval $I \subset \mathbb{R}$ to $\mathbb{R}$ will be denoted by $\partial f$, or $\partial_t f$ if the variable $t$ is well identified. Its value at $t_0 \in I$ is denoted either as $\partial f(t_0)$ or $\partial f|_{t=t_0}$. Higher derivatives are denoted as $\partial^k f$, $k \ge 0$, with the usual convention $\partial^0 f = f$. Note that notation such as $f', f'', f^{(3)}$ will *never* refer to derivatives.

In the following, $U$ is an open subset of $\mathbb{R}^d$. If $f$ is a function from $U$ to $\mathbb{R}^m$, we let $f^{(i)}$ denote the $i^{\text{th}}$ component of $f$, so that

$$f(x) = \begin{pmatrix} f^{(1)}(x) \\ \vdots \\ f^{(m)}(x) \end{pmatrix}$$

for $x \in U$. If $d = 1$, and $f$ is differentiable, the derivative of $f$ at $x$ is the column vector of the derivatives of its components,

$$\partial f(x) = \begin{pmatrix} \partial f^{(1)}(x) \\ \vdots \\ \partial f^{(m)}(x) \end{pmatrix}$$

For $d \ge 1$ and $j \in \{1, \ldots, d\}$, the $j^{\text{th}}$ partial derivative of $f$ at $x$ is

$$\partial_j f(x) = \partial(t \mapsto f(x + t\mathfrak{e}_j))|_{t=0} \in \mathbb{R}^m,$$

where $\mathfrak{e}_1, \ldots, \mathfrak{e}_d$ form the canonical basis of $\mathbb{R}^d$. If the notation for the variables on which $f$ depends is well understood from the context, we will alternatively use $\partial_{x_j} f$. (For example, if $f : (\alpha, \beta) \mapsto f(\alpha, \beta)$, we will prefer $\partial_\alpha f$ to $\partial_1 f$.) The differential of $f$ at $x$ is the linear mapping from $\mathbb{R}^d$ to $\mathbb{R}^m$ represented by the matrix

$$df(x) = [\partial_1 f(x), \ldots, \partial_d f(x)].$$

It is defined so that, for all $h \in \mathbb{R}^d$

$$df(x)h = \partial(t \mapsto f(x + th))|_{t=0}$$

where the right-hand side is the directional derivative of $f$ at $x$ in the direction $h$. Note that, if $f : \mathbb{R}^d \to \mathbb{R}$ (i.e., $m = 1$), $df(x)$ is a row vector. If $f$ is differentiable

on $U$ and $df(x)$ is continuous as a function of $x$, one says that $f$ is continuously differentiable, or $C^1$.

Differentials obey the product rule and the chain rule. If $f, g : U \to \mathbb{R}$, then

$$d(fg)(x) = f(x)dg(x) + g(x)df(x).$$

If $f : U \to \mathbb{R}^m$, $g : \tilde{U} \subset \mathbb{R}^k \to U$, then

$$d(f \circ g)(x) = df(g(x))dg(x).$$

If $d = m$ (so that $df(x)$ is a square matrix), we let $\nabla \cdot f(x) = \text{trace}(df(x))$, the divergence of $f$.

The Euclidean gradient of a differentiable function $f : U \to \mathbb{R}$ is $\nabla f(x) = df(x)^T$. More generally, one defines the gradient of $f$ with respect to a tensor field $x \mapsto A(x)$ taking values in $\mathcal{S}_d^{++}$, as the vector $\nabla_A f(x)$ that satisfies the relation

$$df(x)h = \nabla_A f(x)^T A(x)h$$

for all $h \in \mathbb{R}^d$, so that

$$\nabla_A f(x) = A(x)^{-1} df(x)^T. \tag{1.3}$$

In particular, the Euclidean gradient is associated with $A(x) = \text{Id}_{\mathbb{R}^d}$ for all $x$. With some abuse of notation, we will denote $\nabla_A f = A^{-1} \nabla f$ when $A$ is a fixed matrix, therefore identified with the constant tensor field $x \mapsto A$.

**3.** We here compute, as an illustration and because they will be useful later, the differential of the determinant and the inversion in matrix spaces.

Recall that, if $A = [a_1, \ldots, a_d] \in \mathcal{M}_d$ is a $d$ by $d$ matrix,, with $a_1, \ldots, a_d \in \mathbb{R}^d$, $\det(A)$ is a $d$-linear form $\delta(a_1, \ldots, a_d)$ which vanishes when two columns coincide and such that $\delta(\varepsilon_1, \ldots, \varepsilon_d) = 1$. In particular $\delta$ changes signs when two of its columns are inverted. It follows from this that

$$\partial_{a_{ij}} \det(A) = \delta(a_1, \ldots, a_{i-1}, \varepsilon_j, a_{j+1}, \ldots, a_d)$$
$$= (-1)^{i-1} \delta(\varepsilon_j, a_1, \ldots, a_{i-1}, \ldots, a_d) = (-1)^{i+j} \det A^{(ij)},$$

where $A^{(ij)}$ is the matrix $A$ with row $i$ and column $j$ removed. We therefore find that the differential of $A \mapsto \det(A)$ is the mapping

$$H \mapsto \text{trace}(\text{cof}(A)^T H) \tag{1.4}$$

where $\text{cof}(A)$ is the matrix composed of co-factors $(-1)^{i+j} \det A^{(ij)}$. As a consequence, if $A$ is invertible, then the differential of $\log|\det(A)|$ is the mapping

$$H \mapsto \text{trace}(\det(A)^{-1} \text{cof}(A)^T H) = \text{trace}(A^{-1} H) \tag{1.5}$$

Consider now the function $I(A) = A \mapsto A^{-1}$ defined on $\mathcal{GL}_d(\mathbb{R})$, which is an open subset of $\mathcal{M}_d(\mathbb{R})$. Using $AI(A) = \mathrm{Id}_{\mathbb{R}^d}$ and the product rule, we get

$$A(dI(A)H) + HI(A) = 0$$

or

$$dI(A)H = -A^{-1}HA^{-1}. \tag{1.6}$$

**4.** Higher-order partial derivatives $\partial_{i_k} \cdots \partial_{i_1} f : U \to \mathbb{R}^m$ are defined by iterating the definition of first-order derivatives, namely

$$\partial_{i_k} \cdots \partial_{i_1} f(x) = \partial_{i_k}(\partial_{i_{k-1}} \cdots \partial_{i_1} f)(x)$$

If all order $k$ partial derivatives of $f$ exist and are continuous, one says that $f$ is $k$-times continuously differentiable, or $C^k$ and, when true, the order in which the derivatives are taken does not matter. In this case, one typically groups derivatives with the same order using a power notation, writing, for example

$$\partial_1 \partial_2 \partial_1 f = \partial_1^2 \partial_2 f$$

for a $C^3$ function.

If $f$ is $C^k$, its $k^{\text{th}}$ differential at $x$ is a symmetric $k$-multilinear map that can also be iteratively defined by (for $h_1, \ldots, h_k \in \mathbb{R}^d$)

$$d^k f(x)(h_1, \ldots, h_k) = d(d^{k-1} f(x)(h_1, \ldots, h_{k-1}))h_k \in \mathbb{R}^m.$$

It is related to partial derivatives through the relation:

$$d^k f(x)(h_1, \ldots, h_k) = \sum_{i_1, \ldots, i_k = 1}^{d} h_1^{(i_1)} \cdots h_k^{(i_k)} \partial_{i_k} \cdots \partial_{i_1} f(x).$$

When $m = 1$ and $k = 2$, one denotes by $\nabla^2 f(x) = (\partial_i \partial_j f(x), i, j = 1, \ldots, n)$ the symmetric matrix formed by partial derivatives of order 2 of $f$ at $x$. It is called the *Hessian* of $f$ at $x$ and satisfies

$$h_1^T \nabla^2 f(x) h_2 = d^2 f(x)(h_1, h_2).$$

The Laplacian of $f$ is the trace of $\nabla^2 f$ and denoted $\Delta f$.

**5.** Taylor's theorem, in its integral form, generalizes the fundamental theorem of calculus to higher derivatives. It expresses the fact that, if $f$ is $C^k$ on $U$ and $x, y \in U$ are such that the closed segment $[x, y]$ is included in $U$, then, letting $h = y - x$:

$$f(x + h) = f(x) + df(x)h + \frac{1}{2}d^2 f(x)(h, h) + \cdots + \frac{1}{(k-1)!}d^{k-1} f(x)(h, \ldots, h)$$

$$+ \frac{1}{(k-1)!} \int_0^1 (1-t)^{k-1} d^k f(x + th)(h, \ldots, h)\, dt \tag{1.7}$$

The last term (remainder) can also be written as

$$\frac{1}{k!} \frac{\int_0^1 (1-t)^{k-1} d^k f(x+th)(h,\dots,h)\, dt}{\int_0^1 (1-t)^{k-1}\, dt}.$$

If $f$ takes scalar values, then $d^k f(x+th)(h,\dots,h)$ is real and the intermediate value theorem implies that there exists some $z$ in $[x, y]$ such that

$$f(x+h) = f(x) + df(x)h + \frac{1}{2}d^2 f(x)(h,h) + \cdots + \frac{1}{(k-1)!}d^{k-1} f(x)(h,\dots,h)$$

$$+ \frac{1}{k!}d^k f(z)(h,\dots,h). \quad (1.8)$$

This is not true if $f$ takes vector values. However, for any $M$ such that $|d^k f(z)| \le M$ for $z \in [x, y]$ (such $M$'s always exist because $f$ is $C^k$), one has

$$\frac{1}{(k-1)!} \int_0^1 (1-t)^{k-1} d^k f(x+th)(h,\dots,h)\, dt \le \frac{M}{k!}|h|^k.$$

Equation (1.7) can be written as

$$f(x+h) = f(x) + df(x)h + \frac{1}{2}d^2 f(x)(h,h) + \cdots + \frac{1}{k!}d^k f(x)(h,\dots,h)$$

$$+ \frac{1}{(k-1)!} \int_0^1 (1-t)^{k-1}(d^k f(x+th)(h,\dots,h) - d^k f(x)(h,\cdots,h))\, dt. \quad (1.9)$$

Let

$$\epsilon_x(r) = \max\Big\{|d^k f(x+h) - d^k f(x)| : |h| \le r\Big\}.$$

Since $d^k f$ is continuous, $\epsilon_x(r)$ tends to $0$ when $r \to 0$ and we have

$$\int_0^1 (1-t)^{k-1}|d^k f(x+th)(h,\dots,h) - d^k f(x)(h,\cdots,h)|\, dt \le \frac{|h|^k}{k}\epsilon_x(|h|).$$

This shows that (1.7) implies that

$$f(x+h) = f(x) + df(x)h + \frac{1}{2}d^2 f(x)(h,h) + \cdots + \frac{1}{k!}d^k f(x)(h,\dots,h) + \frac{|h|^k}{k!}\epsilon_x(|h|) \quad (1.10)$$

$$= f(x) + df(x)h + \frac{1}{2}d^2 f(x)(h,h) + \cdots + \frac{1}{k!}d^k f(x)(h,\dots,h) + o(|h|^k) \quad (1.11)$$

## 1.4   Probability theory

**1.** When discussing probabilistic concepts, we will make the convenient assumption that all random variables are defined on a fixed probability space $(\Omega, \mathbb{P})$. This means that $\Omega$ is large enough to include enough randomness to generate all required variables (and implicitly enlarged when needed).

We assume that the reader is familiar with concepts related to discrete random variables or continuous variables (with values in $\mathbb{R}^d$ for some $d$) and their probability density functions, or p.d.f.'s. In particular, $X : \Omega \to \mathbb{R}^d$ is a random variable with p.d.f. $f$ if and only if the expectation of $\varphi(X)$ is given by

$$\mathbb{E}(\varphi(X)) = \int_{\mathbb{R}^d} \varphi(x) f(x) dx$$

for all bounded and continuous functions $\varphi : \mathbb{R}^d \to [0, +\infty)$.

**2.** With a few exceptions, we will use capital letters for random variables and small letters for scalars and vectors that represent realizations of these variables. One of these exceptions will be our notation for training data, defined as an independent and identically distributed (i.i.d.) sample of a given random variable. A realization of such a sample will always be denoted $T = (x_1, \ldots, x_N)$, which is therefore a series of observations. We will use the notation $\mathbb{T} = (X_1, \ldots, X_N)$ for the collection of i.i.d. random variables that generate the training set, so that $T = (X_1(\omega), \ldots, X_N(\omega)) = \mathbb{T}(\omega)$ for some $\omega \in \Omega$. Another exception will apply to variables denoted using Greek letters, for which we will use boldface fonts (such as $\boldsymbol{\alpha}, \boldsymbol{\beta}, \ldots$).

For a random variable $X$, the notation $[X = x]$, or $[X \in A]$ refers to subsets of $\Omega$, for example,

$$[X = x] = \{\omega \in \Omega : X(\omega) = x\}.$$

**3.** As much as possible—but not always—we will avoid making explicit reference to measure theory, leaving to readers familiar with this theory the task to complete the notation and sometimes assumption gaps in order to make some of our statements fully rigorous.

However, there will be situations in which the flexibility of the measure-theoretic formalism is needed for the exposition. The following notions may help the reader navigate through these situations (basic references in measure theory are Rudin [171], Dudley [66], Billingsley [32]).

A measurable space is a pair $(S, \mathcal{S})$ where $S$ is a set and $\mathcal{S} \subset \mathcal{P}(S)$ contains $S$, is stable by complementation (if $A \in \mathcal{S}$, then $A^c = S \setminus A \in \mathcal{S}$), by countable unions and intersections. Such an $\mathcal{S}$ is called a $\sigma$-algebra and elements of $\mathcal{S}$ form the measurable subsets of $S$ (relative to the $\sigma$-algebra).

A (positive) measure $\mu$ on $(S, \mathcal{S})$ in a mapping from $\mathcal{S} \to [0, +\infty)$ that associates to $A \in \mathcal{S}$ its measure $\mu(A)$, such that the measure of a countable union of disjoint sets

is the countable sum of their measures. A function $f : \Omega \to \mathbb{R}^d$ is called measurable if the inverse images by $f$ of open subsets of $\mathbb{R}^d$ are mesurable.

A measurable set $A$ (or event) is negligible (for $\mathbb{P}$) if $\mathbb{P}(A) = 0$ and events are said to happen almost surely if their complements are negligible, i.e., $\mathbb{P}(A^c) = 0$.

**4.** The integral of a function $f : \Omega \to \mathbb{R}^d$ with respect to a measure (such as $\mathbb{P}$) is denoted $\int_S f(x)\mu(dx)$. This integral is defined, using a limit argument, as a function which is linear in $f$ and such that

$$\int_A \mu(dx) = \int_S \mathbf{1}_A(x)\mu(dx) = \mu(A).$$

The Lebesgue measure, $\mathcal{L}_d$, on $\mathbb{R}^d$ provides an important example. For this measure $\mathcal{S}$ is the $\sigma$-algebra generated by open subsets, $\int_{\mathbb{R}^d} f(x)\mathcal{L}_d(dx)$ extends the Riemann integral and is denoted $\int_{\mathbb{R}^d} f(x)dx$. Another important example, when $S$ is finite or countable, is the counting measure, denoted *card*, that return the number of elements of a set, so that $card(A) = |A|$. In this case, $\mathcal{S} = \mathcal{P}(S)$ and the integral is simply the sum:

$$\int_S f(x)card(dx) = \sum_{x \in \mathcal{F}} f(x).$$

**5.** If $\mu$ and $\nu$ are measures on $(S, \mathcal{S})$, one says that $\nu$ is absolutely continuous with respect to $\mu$ and write $\nu \ll \mu$ if,

$$\forall A \in \mathcal{S} : \mu(A) = 0 \Rightarrow \nu(A) = 0. \tag{1.12}$$

The Radon-Nikodym theorem states that $\nu \ll \mu$ if and only if $\nu$ has a density with respect to $\mu$, i.e., there exists a measurable function $\varphi : S \to [0, +\infty)$ such that

$$\int_S f(x)\nu(dx) = \int_S f(x)\varphi(x)\mu(dx)$$

for all measurable $f : S \to [0, +\infty)$.

**6.** If $\mu_1$ is a measure on $(S_1, \mathcal{S}_1)$ and $\mu_2$ a measure on $(S_2, \mathcal{S}_2)$, their tensor product is denoted $\mu_1 \otimes \mu_2$. It is a measure on $S_1 \times S_2$ defined by $\mu_1 \otimes \mu_2(A_1 \times A_2) = \mu_1(A_1)\mu_2(A_2)$ for $A_1 \in \mathcal{S}_1$ and $A_2 \in \mathcal{S}_2$ (the $\sigma$-algebra on $S_1 \times S_2$ is the smallest one that contains all sets $A_1 \times A_2$, $A_1 \in \mathcal{S}_1$, $A_2 \in \mathcal{S}_2$).

The integral, with respect to the product measure, of a function $f : S_1 \times S_2 \to \mathbb{R}^d$ is denoted

$$\int_{S_1 \times S_2} f(x_1, x_2)\mu_1(dx_1)\mu_2(dx_2) = \int_{S_1 \times S_2} f(x_1, x_2)\mu_1 \otimes \mu_2(dx_1, dx_2).$$

The tensor product between more that two measures is defined similarly, with notation

$$\mu_1 \otimes \cdots \otimes \mu_n = \bigotimes_{k=1}^{n} \mu_k.$$

**7.** When using measure-theoretic probability, we will therefore assume that the pair $(\Omega, \mathbb{P})$ is completed to a triple $(\Omega, \mathcal{A}, \mathbb{P})$ where $\mathcal{A}$ is a $\sigma$-algebra and $\mathbb{P}$ a probability measure, that is a positive measure on $(\Omega, \mathcal{A})$ such that $\mathbb{P}(\Omega) = 1$. This triple is called a *probability space*.

A random variable $X$ must then also take values in a measurable space, say $(S, \mathcal{S})$, and must be such that, for all $C \in \mathcal{S}$, the set $[X \in C]$ belongs to $\mathcal{A}$. This justify the computation of $\mathbb{P}(X \in C)$, which will also be denoted $P_X(C)$.

A random variable $X$ taking values in $\mathbb{R}^d$ has a p.d.f. if and only if $P_X \ll \mathcal{L}_d$ and the p.d.f. is the density provided by the Radon-Nikodym theorem. For a discrete random variable (i.e., taking values in a finite or countable set), the p.m.f. of $X$ is also the density of $P_X$ with respect to the counting measure *card*.

If $X$ is a random variable with values in $\mathbb{R}^d$, the integral of $X$ with respect to $\mathbb{P}$ is the expectation of $X$, denoted $\mathbb{E}(X)$. More generally, if $(S, \mathcal{S}, P)$ is a probability space, we will use the notation

$$E_P(f) = \int_S f(x) P(dx).$$

If $P = P_X$ for some random variable $X : \Omega \to S$, we will use $E_X$ rather than $E_{P_X}$.

**8.** One more technical consideration. Whenever we will consider measurable spaces, and sometimes without additional mention, we will assume that these spaces are complete metric spaces that have a dense countable subset (i.e., that are separable). If not specified otherwise, their $\sigma$-algebras are given by the smallest ones containing all open sets (the Borel $\sigma$-algebra).

# Chapter 2

# A Few Results in Matrix Analysis

This chapter collects a few results in linear algebra that will be useful in the rest of this book.

## 2.1 Notation and basic facts

We denote by $\mathcal{M}_{n,d}(\mathbb{R})$ the space of all $n \times d$ matrices with real coefficients[1]. For a matrix $A \in \mathcal{M}_{n,d}(\mathbb{R})$ and integer $k \leq n$ and $l \leq d$, we let $A_{\lceil kl \rceil} \in \mathcal{M}_{k,l}(\mathbb{R})$ denote the matrix $A$ restricted to its first $k$ rows and first $l$ columns. The $i,j$ entry of $A$ will be denoted $A(i,j)$ or $A^{(ij)}$.

We assume that the reader is familiar with elementary matrix analysis, including, in particular the fact that symmetric matrices are diagonalizable in an orthonormal basis, i.e., if $A \in \mathcal{M}_{d,d}(\mathbb{R})$ is a symmetric matrix (whose space is denoted $\mathcal{S}_d$), there exists an orthogonal matrix $U \in \mathcal{O}_d$ (i.e., satisfying $U^T U = U U^T = \mathrm{Id}_{\mathbb{R}^d}$) and a diagonal matrix $D \in \mathcal{M}_{d,d}(\mathbb{R})$ such that

$$A = UDU^T.$$

The identity $AU = UD$ then implies that the columns of $U$ form an orthonormal basis of eigenvectors of $A$.

If $A \in \mathcal{S}_d^+$ is positive semi-definite (i.e., $u^T A u \geq 0$ for all $u \in \mathbb{R}^d$), the entries of $D$ in the decomposition $A = UDU^T$ are non-negative, and one can define the matrix square root of $A$ as $S = UD^{\odot 1/2}U^T$ where $D^{\odot 1/2}$ is the diagonal matrix formed taking the square roots of all coefficients of $D$. We will use the notation $S = A^{1/2}$. Note that $D^{1/2} = D^{\odot 1/2}$ if $D$ is diagonal and positive semi-definite.

If $A \in \mathcal{S}_d^{++}$ is positive definite (i.e., $A$ is positive semi-definite and $u^T A u = 0$ implies $u = 0$) and $B$ is positive semi-definite, both being $d \times d$ matrices, the generalized

---

[1]Unless mentioned otherwise, all matrices are assumed to be real.

eigenvalue problem associated with $A$ and $B$ consists in finding a diagonal matrix $D$ and a matrix $U$ such that $BU = AUD$ and $U^T A U = \mathrm{Id}_{\mathbb{R}^d}$. Letting $\tilde{U} = A^{1/2} U$, the problem is equivalent to solving $A^{-1/2} B A^{-1/2} \tilde{U} = \tilde{U} D$ with $\tilde{U}^T \tilde{U} = \mathrm{Id}_{\mathbb{R}^d}$, i.e., finding the eigenvalue decomposition of the symmetric positive-definite matrix $A^{-1/2} B A^{-1/2}$.

If $A \in \mathcal{M}_{n,d}(\mathbb{R})$, it can be decomposed as

$$A = U D V^T$$

where $U \in \mathcal{O}_n(\mathbb{R})$ and $V \in \mathcal{O}_d(\mathbb{R}))$ are orthogonal matrices and $D \in \mathcal{M}_{n,d}(\mathbb{R})$ is diagonal (i.e., such that $D(i,j) = 0$ whenever $i \neq j$) with non-negative diagonal coefficients. These coefficients are called the singular values of $A$, and the procedure is called a *singular-value decomposition* (SVD) of $A$. An equivalent formulation is that there exist orthonormal bases $u_1, \ldots, u_n$ of $\mathbb{R}^n$ and $v_1, \ldots, v_d$ of $\mathbb{R}^d$ (forming the columns of $U$ and $V$) such that

$$A v_i = \lambda_i u_i$$

for $i \leq \min(n,d)$, where $\lambda_1, \ldots, \lambda_{\min(n,d)}$ are the singular values. Of course, if $A$ is square and symmetric positive semi-definite, an eigenvalue decomposition of $A$ is also a singular value decomposition (and the singular values coincide with the eigenvalues). More generally, if $A = U D V^T$, then $A A^T = U D D^T U^T$ and $A^T A = V D^T D V^T$ are eigenvalue decompositions of $A A^T$ and $A^T A$. Singular values are uniquely defined, up to reordering. However, the matrices $U$ and $V$ are not unique up to column reordering in general.

If $m = \min(n,d)$, then, forming the matrices $\tilde{U} = U_{\lceil n,m \rceil}$ (resp. $\tilde{V} = V_{\lceil d,m \rceil}$) by removing from $U$ (resp. $V$) its last $n - m$ (resp. $d - m$) columns , and $\tilde{D} = D_{\lceil m,m \rceil}$ by removing from $D$ its $n - m$ rows and $d - m$ columns, one has

$$A = \tilde{U} \tilde{D} \tilde{V}^T$$

with $\tilde{U}$, $\tilde{D}$ and $\tilde{V}$ having respectively size $n \times m$, $m \times m$ and $m \times d$, $\tilde{U}^T \tilde{U} = \tilde{V}^T \tilde{V} = \mathrm{Id}_{\mathbb{R}^m}$ and $\tilde{D}$ diagonal with non-negative coefficients. This representation provides a *reduced SVD* of $A$ and one can create a full SVD from a reduced one by completing the missing rows of $\tilde{U}$ and $\tilde{V}$ to form orthogonal matrices, and by adding the required number of zeros to $\tilde{D}$.

## 2.2   The trace inequality

We now descibe Von Neumann's trace theorem. Its justification follows the proof given in Mirsky [137].

**Theorem 2.1 (Von Neumann)** *Let* $A, B \in \mathcal{M}_{n,d}(\mathbb{R})$ *have singular values* $(\lambda_1, \ldots, \lambda_m)$ *and* $(\mu_1, \ldots, \mu_m)$, *respectively, where* $m = \min(n,d)$. *Assume that these eigenvalues are listed*

*in decreasing order so that $\lambda_1 \geq \cdots \geq \lambda_m$ and $\mu_1 \geq \cdots \geq \mu_m$. Then,*

$$\text{trace}(A^T B) \leq \sum_{i=1}^{m} \lambda_i \mu_i. \tag{2.1}$$

*Moreover, if $\text{trace}(A^T B) = \sum_{i=1}^{m} \lambda_i \mu_i$, then there exist $n \times n$ and $d \times d$ orthogonal matrices $U$ and $V$ such that $U^T A V$ and $U^T B V$ are both diagonal, i.e., one can find SVDs of $A$ and $B$ in the same bases of $\mathbb{R}^n$ and $\mathbb{R}^d$.*

PROOF We can assume without loss of generality that $d \leq n$ because, if the result holds for $A$ and $B$, it also holds for $A^T$ and $B^T$. Let $A = U_1 \Lambda V_1^T$ and $B = U_2 M V_2^T$ be the singular values decompositions of $A$ and $B$ (both $\Lambda$ and $M$ are $n \times d$ matrices). Then

$$\text{trace}(A^T B) = \text{trace}(V_1 \Lambda^T U_1^T U_2 M V_2) = \text{trace}(\Lambda^T U M V^T)$$

with $U = U_1^T U_2$ and $V = V_1^T V_2$. Let $u(i,j), 1 \leq i,j \leq n$ and $v(i,j), 1 \leq i,j \leq d$ be the coefficients of the orthogonal matrices $U$ and $V$. Then

$$\text{trace}(\Lambda^T U M V^T) = \sum_{i,j=1}^{d} u(i,j)v(i,j)\lambda_i \mu_j \leq \frac{1}{2}\sum_{i,j=1}^{d} \lambda_i \mu_j u(i,j)^2 + \frac{1}{2}\sum_{i,j=1}^{d} \lambda_i \mu_j v(i,j)^2 \tag{2.2}$$

Let us consider the first sum in the upper-bound. Let $\xi_d = \lambda_d$ (resp. $\eta_d = \mu_d$) and $\xi_i = \lambda_i - \lambda_{i+1}$ (resp. $\eta_i = \mu_i - \mu_{i+1}$) for $i = 1,\ldots,d-1$. Since singular values are non-increasing, we have $\xi_i, \eta_i \geq 0$ and

$$\lambda_i = \sum_{j=i}^{d} \xi_j, \quad \mu_i = \sum_{j=i}^{d} \eta_j$$

for $i = 1,\ldots,d$. We have

$$\sum_{i,j=1}^{d} \lambda_i \mu_j u(i,j)^2 = \sum_{i,j=1}^{d} \sum_{i'=i}^{d} \xi_{i'} \sum_{j'=j}^{d} \eta_{j'} u(i,j)^2 = \sum_{i',j'=1}^{d} \xi_{i'} \eta_{j'} \sum_{i=1}^{i'} \sum_{j=1}^{j'} u(i,j)^2$$

$$\leq \sum_{i',j'=1}^{d} \xi_{i'} \eta_{j'} \min(i',j') \tag{2.3}$$

where we used the fact that $U$ is orthogonal, which implies that $\sum_{j=1}^{j'} u(i,j)^2$ and $\sum_{i=1}^{i'} u(i,j)^2$ are both less than 1. Notice also that, when $u(i,j) = \delta_{ij}$ (i.e., $u(i,j) = 1$ if $i = j$ and zero otherwise), then

$$\sum_{i=1}^{i'} \sum_{j=1}^{j'} u(i,j)^2 = \min(i',j'),$$

so that the last inequality is an identity, and the chain of equalities leading to (2.3) implies

$$\sum_{i',j'=1}^{d} \xi_{i'} \eta_{j'} \min(i',j') = \sum_{i=1}^{d} \lambda_i \mu_j.$$

We therefore obtain (for any $U$), the fact that

$$\sum_{i,j=1}^{d} \lambda_i \mu_j u(i,j)^2 \leq \sum_{i=1}^{d} \lambda_i \mu_j.$$

The same identity obviously holds with $v$ in place of $u$, and combining the two yields (2.1).

We now consider conditions for equality. Clearly, if one can find SVD decompositions of $A$ and $B$ with $U_1 = U_2$ and $V_1 = V_2$, then $U = \mathrm{Id}_{\mathbb{R}^n}$, $V = \mathrm{Id}_{\mathbb{R}^d}$ and (2.1) is an identity. We want to prove the converse statement.

For (2.1) to be an equality, we first need (2.2) to be an identity, which requires that $u(i,j) = v(i,j)$ as soon as $\lambda_i \mu_j > 0$. We also need an equality in (2.3), which requires

$$\sum_{i=1}^{i'} \sum_{j=1}^{j'} u(i,j)^2 = \min(i',j')$$

as soon as $\lambda_{i'} > \lambda_{i'+1}$ and $\mu_{j'} > \mu_{j'+1}$. The same identity must be true with $v(i,j)$ replacing $u(i,j)$

In view of this, denote by $i_1 < \cdots < i_p$ (resp. $j_1 < \cdots < j_q$) the indexes at which the singular values of $A$ (resp. $B$) differ form their successors, with the convention $\lambda_{d+1} = \mu_{d+1} = 0$. Let, for $k = 1,\ldots,p$ and $l = 1,\ldots,q$

$$C(k,l) = \sum_{i=1}^{i_k} \sum_{j=1}^{j_l} u(i,j)^2.$$

Then, we must have $C(k,l) = \min(i_k,j_l)$ for all $k,l$ and $u(i,j) = v(i,j)$ for $i = 1,\ldots,i_p$ and $j = 1,\ldots,j_q$.

If, for all $i,j \leq d$, we let $U_{\lceil ij \rceil}$ be the matrix formed by the first $i$ rows and $j$ columns of $U$, the condition $C_{kl} = \min(i_k,j_l)$ requires that $U_{\lceil i_k j_l \rceil} U_{\lceil i_k j_l \rceil}^T = \mathrm{Id}_{\mathbb{R}^{i_k}}$ if $i_k \leq j_l$ and $U_{\lceil i_k j_l \rceil}^T U_{\lceil i_k j_l \rceil} = \mathrm{Id}_{\mathbb{R}^{j_l}}$ if $j_l \leq i_k$. This shows that, if $i_k \leq j_l$, the rows of $U_{\lceil i_k j_l \rceil}$ form an orthonormal family, and necessarily, all elements $u(i,j)$ for $i \leq i_k$ and $j > j_l$ vanish. The symmetric situation holds if $j_l \leq i_k$.

Let $r_k = i_k - i_{k-1}$ and $s_l = j_l - j_{l-1}$ (with $i_0 = j_0 = 0$). We now consider possible changes in the SVDs of $A$ and $B$. With our notation, the matrix $\Lambda$ takes the form

$$
\Lambda = \begin{pmatrix}
\lambda_{i_1} \mathrm{Id}_{\mathbb{R}^{r_1}} & 0 & 0 & \cdots & 0 & 0 & \ldots & 0 \\
0 & \lambda_{i_2} \mathrm{Id}_{\mathbb{R}^{r_2}} & 0 & \cdots & 0 & 0 & \ldots & 0 \\
\vdots & & \ddots & & \vdots & \vdots & & \vdots \\
0 & 0 & \ldots & \lambda_{i_p} \mathrm{Id}_{\mathbb{R}^{r_p}} & 0 & 0 & \ldots & 0 \\
0 & & \cdots & & 0 & 0 & \ldots & 0 \\
\vdots & & & & \vdots & \vdots & & \vdots \\
0 & & \cdots & & 0 & 0 & \ldots & 0
\end{pmatrix}
$$

Let $W, \tilde{W}$ be $n \times n$ and $d \times d$ orthogonal matrices taking the form

$$
W = \begin{pmatrix}
W_1 & 0 & 0 & \cdots & 0 \\
0 & W_2 & 0 & \cdots & 0 \\
\vdots & & \ddots & & \vdots \\
0 & 0 & \ldots & W_p & 0 \\
0 & & \cdots & & W_{p+1}
\end{pmatrix}, \quad
\tilde{W} = \begin{pmatrix}
W_1 & 0 & 0 & \cdots & 0 \\
0 & W_2 & 0 & \cdots & 0 \\
\vdots & & \ddots & & \vdots \\
0 & 0 & \ldots & W_p & 0 \\
0 & & \cdots & & \tilde{W}_{p+1}
\end{pmatrix}
$$

where $W_1, \ldots, W_p$ are orthogonal with respective sizes $r_1, \ldots, r_p$, $W_{p+1}$ is orthogonal with size $n - i_p$ and $\tilde{W}_{p+1}$ is orthogonal with size $d - i_p$. Then we have

$$
WD\tilde{W} = D
$$

proving that $U_1$ can be replaced by $U_1 W$ provided that $V_1$ is replaced by $V_1 \tilde{W}$. Similar transformations can be made on $U_1$ and $V_2$, with $U_2$ replaced by $U_2 Z$ and $V_2$ by $V_2 \tilde{Z}$ with

$$
Z = \begin{pmatrix}
Z_1 & 0 & 0 & \cdots & 0 \\
0 & Z_2 & 0 & \cdots & 0 \\
\vdots & & \ddots & & \vdots \\
0 & 0 & \ldots & Z_q & 0 \\
0 & & \cdots & & Z_{q+1}
\end{pmatrix}, \quad
\tilde{Z} = \begin{pmatrix}
Z_1 & 0 & 0 & \cdots & 0 \\
0 & Z_2 & 0 & \cdots & 0 \\
\vdots & & \ddots & & \vdots \\
0 & 0 & \ldots & Z_q & 0 \\
0 & & \cdots & & \tilde{Z}_{q+1}
\end{pmatrix}
$$

with a structure similar to $W$ and $\tilde{W}$, replacing $r_1, \ldots, r_p$ by $s_1, \ldots, s_q$. As a consequence, $U = U_1^T U_2$ can be replaced by $W^T U Z$ and $V$ by $\tilde{W}^T V \tilde{Z}$. To complete the proof, we need to show that, when (2.1) is an equality, these matrices can be chosen so that $W^T U Z = \mathrm{Id}_{\mathbb{R}^n}$ and $\tilde{W}^T V \tilde{Z} = \mathrm{Id}_{\mathbb{R}^d}$.

Let us consider a first step in this direction, assuming that $i_1 \leq j_1$ so that

$$
U_{[i_1 j_1]} U_{\lceil i_1 j_1 \rceil}^T = \mathrm{Id}_{\mathbb{R}^{i_1}}.
$$

Complete $U^T_{\lceil i_1 j_1 \rceil}$ into a orthogonal matrix $Z_1 = [U^T_{\lceil i_1 j_1 \rceil}, \tilde{U}]$. Build a matrix $Z$ as above by taking $Z_2, \ldots, Z_{q+1}$ equal to the identity. Then $UZ$ has a first $i_1 \times i_1$ block equal to $\mathrm{Id}_{\mathbb{R}^{i_1}}$, which implies that all coefficients on the right and below this block are zeros. If $j_1 \leq i_1$, a similar construction can be made on the other side, letting $W_1 = [U_{\lceil i_1 j_1 \rceil} \tilde{U}]$ with the first $j_1 \times j_1$ block of the new matrix $U$ equal to the identity. Note that, since $V_{\lceil i_p j_q \rceil} = U_{\lceil i_p j_q \rceil}$, the same result is obtained on $V$ at the same time.

Pursuing this way (and skipping the formal induction argument, which is a bit tedious), we can progressively introduce identity blocks into $U$ and $V$ and transform them into new matrices (that we still denote by $U$ and $V$) taking the form (letting $k = \min(i_p, j_q)$)

$$U = \begin{pmatrix} \mathrm{Id}_{\mathbb{R}^k} & 0 \\ 0 & \bar{U} \end{pmatrix} \text{ and } V = \begin{pmatrix} \mathrm{Id}_{\mathbb{R}^k} & 0 \\ 0 & \bar{V} \end{pmatrix}$$

If $k = i_p$ (resp. $k = j_q$), the final reduction can be obtained by choosing $W_{p+1} = \bar{U}$ and $\tilde{W}_{p+1} = \bar{V}$ (resp. $Z_{p+1} = \bar{U}^T$ and $\tilde{Z}_{p+1} = \bar{V}^T$), leading to SVDs for $A$ and $B$ with identical matrices $U_1 = U_2$ and $V_1 = V_2$.                                   ∎

**Remark 2.2** Note that, since the singular values of $-A$ and of $A$ coincide, theorem 2.1 implies

$$\left| \mathrm{trace}(A^T B) \right| \leq \sum_{i=1}^m \lambda_i \mu_i. \tag{2.4}$$

for all matrices $A$ and $B$, with equality if either $A$ and $B$ or $-A$ and $B$ have an SVD using the same bases.                                                                    ♦

## 2.3  Applications

Let $p$ and $d$ be integers with $p \leq d$. Let $A \in \mathcal{S}_d(\mathbb{R})$, $B \in \mathcal{S}_p(\mathbb{R})$ be symmetric matrices. We consider the following optimization problem: maximize, over matrices $U \in \mathcal{M}_{d,p}(\mathbb{R})$ such that $U^T U = \mathrm{Id}_{\mathbb{R}^p}$, the function

$$F(U) = \mathrm{trace}(U^T A U B) = \mathrm{trace}(A U B U^T).$$

We first note that the singular values of $U B U^T$, which is $d \times d$, are the same as the eigenvalues of $B$ completed with zeros. Letting $\lambda_1 \geq \cdots \geq \lambda_d$ be the eigenvalues of $A$ and $\mu_1 \geq \cdots \geq \mu_p$ those of $B$, we therefore have, from theorem 2.1,

$$F(U) \leq \sum_{i=1}^p \lambda_i \mu_i.$$

Introduce the eigenvalue decompositions of $A$ and $B$ in the form $A = V\Lambda V^T$ and $B = WMW^T$. For $F(U)$ to be equal to its upper-bound, we know that we must arrange $UBU^T$ to take the form

$$UBU^T = V\begin{pmatrix} M & 0 \\ 0 & 0 \end{pmatrix}V^T.$$

Use, as before, the notation $V_{\lceil dp \rceil}$ to denote the matrix formed with the $p$ first columns of $V$. Take $U = V_{\lceil dp \rceil}W^T$, which satisfies $U^T U = \mathrm{Id}_{\mathbb{R}^p}$. We then have

$$V_{\lceil dp \rceil}W^T BW V_{\lceil dp \rceil}^T = V_{\lceil dp \rceil}M V_{\lceil dp \rceil}^T = V\begin{pmatrix} M & 0 \\ 0 & 0. \end{pmatrix}V^T,$$

which shows that $U$ is optimal. We summarize this discussion in the next theorem.

**Theorem 2.3** *Let $A \in \mathcal{S}_d(\mathbb{R})$ and $B \in \mathcal{S}_p(\mathbb{R})$ be symmetric matrices, with $p \le d$. Let eigenvalue decompositions of $A$ and $B$ be given by $A = V\Lambda V^T$ and $B = WMW^T$, where the diagonal elements of $\Lambda$ (resp. $M$) are $\lambda_1 \ge \cdots \ge \lambda_d$ (resp. $\mu_1 \ge \cdots \ge \mu_p$).*

*Define $F(U) = \mathrm{trace}(AUBU^T)$, for $U \in \mathcal{M}_{d,p}(\mathbb{R})$. Then,*

$$\max\left\{F(U): U^T U = \mathrm{Id}_{\mathbb{R}^p}\right\} = \sum_{i=1}^{p} \lambda_i \mu_i.$$

*This maximum is attained at*

$$U = V_{\lceil d,p \rceil}W^T.$$

The following corollary applies theorem 2.3 with $B = \mathrm{diag}(\mu_1, \ldots, \mu_p)$.

**Corollary 2.4** *Let $A \in \mathcal{S}_d(\mathbb{R})$ be a symmetric matrix with eigenvalues $\lambda_1 \ge \cdots \ge \lambda_d$. For $p \le d$, let $\mu_1 \ge \cdots \ge \mu_p > 0$ and define*

$$F(e_1, \ldots, e_p) = \sum_{i=1}^{p} \mu_i e_i^T A e_i.$$

*Then, the maximum of F over all orthonormal families $e_1, \ldots, e_p$ in $\mathbb{R}^d$ is $\sum_{i=1}^{p} \lambda_i \mu_i$ and is attained when $e_1, \ldots, e_p$ are eigenvectors of A with eigenvalues $\lambda_1, \ldots, \lambda_p$.*

*The minimum of F over all orthonormal families $e_1, \ldots, e_p$ in $\mathbb{R}^d$ is $\sum_{i=1}^{p} \lambda_{d-i+1} \mu_i$ and is attained when $e_1, \ldots, e_p$ are eigenvectors of A with eigenvalues $\lambda_d, \ldots, \lambda_{d-p+1}$.*

PROOF The statement about the maximum is just a special case of theorem 2.3, with $B = \mathrm{diag}(\mu_1, \ldots, \mu_p)$, noting that the $i$th diagonal element of $U^T A U$ is $e_i^T A e_i$ where $e_i$ is the $i$th column of $U$.

The statement about the minimum is deduced by replacing $A$ by $-A$. ∎

Applying this corollary with $p = 1$, we retrieve the elementary result that $\lambda_1 = \max\{u^T A u : |u| = 1\}$ and $\lambda_d = \min\{u^T A u : |u| = 1\}$.

To complete this chapter, we quickly state and prove Rayleigh's theorem.

**Theorem 2.5** *Let $A \in \mathcal{M}_{d,d}(\mathbb{R})$ be a symmetric matrix with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_d$. Then*

$$\lambda_k = \max_{V:\dim(V)=k} \min\{u^T A u, u \in V, |u| = 1\} = \min_{V:\dim(V)=d-k+1} \max\{u^T A u, u \in V, |u| = 1\}$$

*where the min and max are taken over linear subspaces of $\mathbb{R}^d$.*

Proof Let $e_1, \ldots, e_d$ be an orthonormal basis of eigenvectors of $A$ associated with $\lambda_1, \ldots, \lambda_d$. Let, for $k \leq l$, $W_{k,l} = \text{span}(e_k, \ldots, e_l)$. Let $V$ be a subspace of dimension $k$. Then $V \cap W_{k,d} \neq \emptyset$ (because the sum of the dimensions of these two spaces is $d + 1$). Taking $u_0$ with norm 1 in this intersection, we have

$$\min\{u^T A u, u \in V, |u| = 1\} \leq u_0^T A u_0 \leq \max\{u^T A u, u \in W_{k,d}, |u| = 1\} = \lambda_k,$$

where the last identity follows by considering the eigenvalues of $A$ restricted to $W_{k,d}$. So, the maximum of the right-hand side is indeed less than $\lambda_k$, and it is attained for $V = W_{1,k}$. This proves the first identity, and the second one can be obtained by applying the first one to $-A$. ∎

## 2.4   Some matrix norms

The *operator norm* of a matrix $A \in \mathcal{M}_{n,d}(\mathbb{R})$, is defined as

$$|A|_{\text{op}} = \max\{|Ax| : x \in \mathbb{R}^d, |x| = 1\}.$$

It is equal to the square root of the largest eigenvalue of $A^T A$, i.e., to the largest singular value of $A$.

The *Frobenius norm* of $A$ is

$$|A|_F = \sqrt{\text{trace}(A^T A)} = \sqrt{\sum_{i,j=1}^{d} A(i,j)^2},$$

so that

$$|A|_F = \left(\sum_{k=1}^{m} \sigma_k^2\right)^{1/2}$$

where $\sigma_1,\ldots,\sigma_m$ are the singular values of $A$ (and $m = \min(n,d)$).

The *nuclear norm* of $A$ is defined by

$$|A|_* = \sum_{k=1}^{d} \sigma_k.$$

One can prove that this is a norm using an equivalent definition, provided by the following proposition.

**Proposition 2.6** *Let $A$ be an $n$ by $d$ matrix. Then*

$$|A|_* = \max\left\{\mathrm{trace}(UAV^T) : U \in \mathcal{M}_{n,n} \text{ and } U^TU = \mathrm{Id}, V \in \mathcal{M}_{d,d} \text{ and } V^TV = \mathrm{Id}\right\}.$$

PROOF The fact that $\mathrm{trace}(UAV^T) \leq |A|_*$ for any $U$ and $V$ is a consequence of the trace inequality applied with $B = [\mathrm{Id}, 0]$ or its transpose depending on whether $n \leq d$ or not. The upper-bound being attained when $U$ and $V$ are the matrices forming the singular value decomposition of $A$, the proof is complete. ∎

The fact that $|A|_*$ is a norm, for which the only non-trivial fact was the triangular inequality, now is an easy consequence of this proposition, because the maximum of the sum of two functions is always less than the sum of their maximums. More precisely, we have

$$
\begin{aligned}
|A + B|_* &= \max\{\mathrm{trace}(UAV^T) + \mathrm{trace}(UBV^T) : \\
&\qquad U^TU = \mathrm{Id}, V^TV = \mathrm{Id}\} \\
&\leq \max\{\mathrm{trace}(UAV^T) : U^TU = \mathrm{Id}, V^TV = \mathrm{Id}\} \\
&\qquad + \max\{\mathrm{trace}(UBV^T) : U^TU = \mathrm{Id}, V^TV = \mathrm{Id}\} \\
&= |A|_* + |B|_*
\end{aligned}
$$

The nuclear norms is also called the Ky Fan norm of order $d$. Ky Fan norms of order $k$ (for $1 \leq k \leq d$) associate to a matrix $A$ the quantity

$$|A|_{(k)} = \lambda_1 + \cdots + \lambda_k,$$

i.e., the sum of its $k$ largest singular values. One has the following proposition.

**Proposition 2.7** *The Ky Fan norms satisfy the triangular inequality.*

PROOF We prove this following the argument suggested in Bhatia [28]. For $A \in \mathcal{M}_{d,d}$, and $k = 1,\ldots,d$, let $\mathrm{trace}_{(k)}(A)$ be the sum of the $k$ largest diagonal elements of $A$. Let,

for a symmetric matrix $A$, $|A|'_{(k)}$ denote the sum of the $k$ largest eigenvalues of $A$ (it is equal to $|A|_{(k)}$ if $A$ is positive definite, but can also include negative values).

Then, for any symmetric matrix $A \in \mathcal{S}_d$,

$$|A|'_{(k)} = \max\left\{\mathrm{trace}_{(k)}(UAU^T) : U \in CO_d\right\}. \tag{2.5}$$

To show this, assume that $V$ in $\mathcal{O}_d$ diagonalizes $A$, so that $D = VAV^T$ is a diagonal matrix. Assume, without loss of generality, that the coefficients $\lambda_j = D(j,j)$ are non-increasing. Fix $U \in \mathcal{O}_d$, let $B = UAU^T$ and $W = VU^T$ so that $D = WBW^T$, or $B = W^TDW$. Then, for any $j \leq d$,

$$B(j,j) = \sum_{i=1}^{d} W(i,j)^2 D(i,i).$$

Then, for any $1 \leq j_1 < \cdots < j_k \leq d$

$$\sum_{l=1}^{k} B(j_l, j_l) = \sum_{i=1}^{d} D(i,i) \sum_{l=1}^{k} W(i,j_l)^2$$

$$= \sum_{i=1}^{k} D(i,i) + \sum_{i=1}^{k} D(i,i)\left(\sum_{l=1}^{k} W(i,j_l)^2 - 1\right) + \sum_{i=k+1}^{d} D(i,i)\sum_{l=1}^{k} W(i,j_l)^2$$

$$= \sum_{i=1}^{k} D(i,i) + \sum_{i=1}^{k}(D(i,i) - D(k,k))\left(\sum_{l=1}^{k} W(i,j_l)^2 - 1\right)$$

$$+ \sum_{i=k+1}^{d}(D(i,i) - D(k,k))\sum_{l=1}^{k} W(i,j_l)^2 + D(k,k)\left(\sum_{i=1}^{n}\sum_{j=1}^{k} W(i,j_l)^2 - k\right).$$

Because $W$ is orthogonal, we have $\sum_{l=1}^{k} W(i,j_l)^2 \leq 1$ and

$$\sum_{i=1}^{n}\sum_{j=1}^{k} W(i,j_l)^2 = k.$$

This shows that the terms after $\sum_{i=1}^{k} D(i,i)$ in the upper bound are negative or zero, so that

$$\sum_{l=1}^{k} B(j_l, j_l) \leq \sum_{i=1}^{k} D(i,i).$$

The maximum of the left-hand side is $\mathrm{trace}_{(k)}(B)$. Noting that we get an equality when choosing $U = V$, the proof of (2.5) is complete.

Using the same argument as that made above for the nuclear norm, one deduces from this that

$$|A + B|'_{(k)} \le |A|'_{(k)} + |B|'_{(k)}$$

for all $A, B \in \mathcal{S}_d$ and all $k = 1, \dots, d$.

Now, let $A \in \mathcal{M}_{n,d}$ and consider the symmetric matrix

$$\tilde{A} = \begin{pmatrix} 0 & A^T \\ A & 0 \end{pmatrix} \in \mathcal{S}_{n+d}.$$

Write a vector $u \in \mathbb{R}^{n+d}$ as $u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$ with $u_1 \in \mathbb{R}^d$ and $u_2 \in \mathbb{R}^n$. Then $u$ is an eigenvector of $\tilde{A}$ for an eigenvalue $\lambda$ if and only if $A^T u_2 = \lambda u_1$ and $A u_1 = \lambda u_2$, which implies that $A^T A u_1 = \lambda^2 u_1$ and $\lambda^2$ is a singular value of $A$. Conversely, if $\mu$ is a nonzero singular value of $A$, associated with eigenvector $u_1$, then $1/\sqrt{\mu}$ and $-1/\sqrt{\mu}$ are eigenvalues of $\tilde{A}$, associated with eigenvectors $\begin{pmatrix} u_1 \\ \pm A u_1/\sqrt{\mu} \end{pmatrix}$. It follows from this that $|A|_{(k)} = |\tilde{A}|'_{(k)}$ for $k \le \min(n, d)$ and therefore satisfies the triangle inequality. ∎

We refer to [28] for more examples of matrix norms, including, in particular those provided by taking $p$th powers in Ky Fan's norms, defining

$$|A|_{(k,p)} = (\lambda_1^p + \cdots + \lambda_k^p)^{1/p}.$$

# Chapter 3

# Introduction to Optimization

This chapter summarizes some fundamental concepts in optimization that will be used later in the book. The reader is referred to textbooks, such as Beck [22], Eiselt et al. [68], Nocedal and Wright [146], Boyd et al. [40] and many others for proofs and deeper results.

## 3.1 Basic Terminology

**1.** If $I$ is a subset of $\mathbb{R}$, a lower bound of $I$ is an element $u \in [-\infty, +\infty]$ such that $u \le x$ for all $x \in I$. Among these lower bounds, there exists a largest element, denoted $\inf I \in [-\infty, +\infty]$, called the infimum of $I$ (by convention, the infimum of an empty set is $+\infty$). Similarly, one defines the supremum of $I$, denoted $\sup I$, as the smallest upper bound of $I$ (and the supremum of an empty set is $-\infty$). Every set in $\mathbb{R}$ has an infimum and a supremum, but these numbers do not necessarily belong to $I$. When they do, they are respectively called minimal and maximal elements of $I$, and are denoted $\min I$ and $\max I$. So, the statement "$u = \min I$" means $u \in I$ and $u \le v$ for all $v \in I$.

**2.** If $F : \Omega \to \mathbb{R}$ is a real-valued function defined on a subset $\Omega \subset \mathbb{R}^d$, the infimum of $F$ over $\Omega$ is defined by

$$\inf_{\Omega} F = \inf\{F(x) : x \in \Omega\}$$

and its supremum is

$$\sup_{\Omega} F = \sup\{F(x) : x \in \Omega\}.$$

As seen above both numbers are well defined, and can take infinite values. One says that $x \in \Omega$ is a (global) minimizer (resp. maximizer) of $F$ if $F(y) \ge F(x)$ (resp. $F(y) \le F(x)$) for all $y \in \Omega$. One also says that $F$ reaches its minimum (resp. maximum), or is minimized (resp. maximized) at $x$. Equivalently, $x$ is a minimizer (resp. maximizer) of $F$ if and only if $x \in \Omega$ and

$$F(x) = \min\{F(y) : y \in \Omega\} \text{ (resp. } \max\{F(y) : y \in \Omega\}\text{)}.$$

In such cases, one also writes $F(x) = \min_\Omega F$ or $F(x) = \max_\Omega F$. In particular, the notation $u = \min_\Omega F$ indicates that $u = \inf_\Omega F$ and that there exists an $x$ in $\Omega$ such that $F(x) = u$ (i.e., that the infimum of $F$ over $\Omega$ is realized at some $x \in \Omega$). Note that the infimum of a function always exists, but not necessarily its minimum. Also note that minimizers, when they exist, are not necessarily unique. We will denote by $\operatorname{argmin}_\Omega F$ (resp. $\operatorname{argmax}_\Omega F$) the (possibly empty) set of minimizers (resp. maximizers) of $F$

**3.** One says that $x$ is a local minimizer (resp. maximizer) of $F$ on $\Omega$ if there exists an open ball $B \subset \mathbb{R}^d$ such that $x \in B$ and $F(x) = \min_{\Omega \cap B} F$ (resp. $F(x) = \max_{\Omega \cap B} F$).

**4.** An optimization problem consists in finding a minimizer or maximizer of an "objective function" $F$. Focusing from now on on minimization problems (statements for maximization problems are symmetric), we will always implicitly assume that a minimizer exists. The following provides some general assumptions on $F$ and $\Omega$ that ensure this fact.

The sublevel sets of $F$ in $\Omega$ are denoted $[F \le u]_\Omega$ (or simply $[F \le u]$ when $\Omega = \mathbb{R}^d$) for $u \in [-\infty, +\infty]$ with

$$[F \le u]_\Omega = \{x \in \Omega : F(x) \le u\}.$$

Note that

$$\operatorname*{argmin}_\Omega F = \bigcap_{u > \inf F} [F \le u]_\Omega.$$

A typical requirement for $F$ is that its sublevel sets are closed in $\mathbb{R}^d$, which means that, if a sequence $(x_n)$ in $\Omega$ satisfies, for some $u \in \mathbb{R}$, $F(x_n) \le u$ for all $n$ and converges to a limit $x$, then $x \in \Omega$ and $F(x) \le u$. If this is true, one says that $F$ is *lower semi-continuous*, or l.s.c, on $\Omega$. If, in addition to being closed, the sublevel sets of $F$ are bounded (at least for $u$ small enough—larger than $\inf F$), then $\operatorname{argmin}_\Omega F$ is an intersection of nested compact sets, and is therefore not empty (so that the optimization problem has at least one solution).

**5.** Different assumptions on $F$ and $\Omega$ lead to different types of minimization problems, with specific underlying theory and algorithms.

  1. If $F$ is $C^1$ or smoother and $\Omega = \mathbb{R}^d$, one speaks of an unconstrained smooth optimization problem.

  2. For constrained problems, $\Omega$ is often specified by a finite number of inequalities, i.e.,

$$\Omega = \{x \in \mathbb{R}^d : \gamma_i(x) \le 0, i = 1, \dots, q\}.$$

If $F$ and all functions $\gamma_1, \dots, \gamma_q$ are $C^1$ one speaks of smooth constrained problems.

  3. If $\Omega$ is a convex set (i.e., $x, y \in \Omega \Rightarrow [x, y] \in \Omega$, where $[x, y]$ is the closed line segment connecting $x$ and $y$) and $F$ is a convex function (i.e., $F((1-t)x + ty) \le (1-t)F(x) + tF(y)$ for all $x, y \in \Omega$), one speaks of a convex optimization problem.

4. Non-smooth problems are often considered in data science, and lead to interesting algorithms and solutions.

5. When both $F$ and $\gamma_1, \dots, \gamma_q$ are affine functions, one speaks of a linear programming problem (or a linear program). (An affine function is a mapping $x \mapsto b^T x + \beta$, $b \in \mathbb{R}^d$, $\beta \in \mathbb{R}$.)

If $F$ is quadratic ($F(x) = \frac{1}{2}x^T A x - b^T x$), and all $\gamma_i$'s are affine, one speaks of a quadratic programming problem.

6. Finally, some machine learning problems are specified over discrete or finite sets $\Omega$ (for example $\mathbb{Z}^d$, or $\{0,1\}^d$), leading to combinatorial optimization problems.

## 3.2 Unconstrained Optimization Problems

### 3.2.1 Conditions for optimality (general case)

Consider a function $F : \Omega \to \mathbb{R}$ where $\Omega$ is an *open* subset of $\mathbb{R}^d$. We first discuss the unconstrained optimization problem of finding

$$x^* \in \underset{\Omega}{\operatorname{argmin}} F. \tag{3.1}$$

The following result summarizes (non-identical) necessary and sufficient conditions that are applicable to such a solution.

**Theorem 3.1 Necessary conditions.** *Assume that $F$ is differentiable over $\Omega$, and that $x^*$ is a local minimum of $F$. Then $\nabla F(x^*) = 0$.*

*If $F$ is $C^2$, then, in addition, $\nabla^2 F(x^*)$ must be positive semidefinite.*

**Sufficient conditions.** *Assume that $F \in C^2(\Omega)$. If $x^* \in \Omega$ is such that $\nabla F(x^*) = 0$ and $\nabla^2 F(x^*)$ is positive definite, then $x^*$ is a local minimum of $F$.*

PROOF Necessary conditions: Since $\Omega$ is open, it contains an open ball centered at $x^*$, with radius $\epsilon_0$ and therefore all segments $[x^*, x^* + \epsilon h]$ for all $\epsilon \in [0, \epsilon_0]$ and all unit norm vectors $h$. Since $x^*$ is a local minimum, we can choose $\epsilon_0$ so that $F(x^* + \epsilon h) \geq F(x^*)$ for all $h$ with $|h| = 1$.

Using Taylor formula, we get (for $\epsilon \in [0, \epsilon_0]$, $|h| = 1$)

$$0 \leq F(x^* + \epsilon h) - f(x^*) = \epsilon \int_0^1 dF(x^* + t\epsilon h) h \, dt.$$

If $dF(x^*)h \neq 0$ for some $h$, then, for small enough $\epsilon$, $dF(x^* + t\epsilon h)h$ cannot change sign for $t \in [0,1]$ and therefore $\int_0^1 dF(x^* + t\epsilon h)h \, dt$ has the same sign as $dF(x^*)(h)$ which must therefore be positive. But the same argument can be made with $h$ replaced by

$-h$, implying that $dF(x^*)(-h) = -dF(x^*)h$ is also positive, and this gives a contradiction. We therefore have $dF(x^*)(h) = 0$ for all $h$, i.e., $\nabla F(x^*) = 0$.

Assume that $F$ is $C^2$. Then, making a second-order Taylor expansion, one gets

$$0 \leq F(x^* + \epsilon h) - F(x^*) = \epsilon^2 \int_0^1 (1-t)d^2F(x^* + t\epsilon h)(h,h)dt.$$

The same argument as above shows that, if $d^2F(x^*)(h,h) \neq 0$, then it must be positive. This shows that $d^2F(x^*)(h,h) \geq 0$ for all $h$ and $d^2F(x^*)$ (or its associated matrix $\nabla^2 F(x^*)$) is positive semidefinite.

Now, assume that $F$ is $C^2$ and $\nabla^2 F(x^*)$ positive definite. One still has

$$F(x^* + \epsilon h) - F(x^*) = \epsilon^2 \int_0^1 (1-t)d^2F(x^* + t\epsilon h)(h,h)dt$$

If $\nabla^2 F(x^*) > 0$, then $\nabla^2 F(x^* + t\epsilon h) > 0$ for small enough $\epsilon$, showing the the r.h.s. of the identity is positive for $h \neq 0$, and that $F(x^* + \epsilon h) > F(x^*)$. ∎

Because maximizing $F$ is the same as minimizing $-F$, necessary (resp. sufficient) conditions for optimality in maximization problems are immediately deduced from the above: it suffices to replace positive semidefinite (resp. positive definite) by negative semidefinite (resp. negative definite).

### 3.2.2   Convex sets and functions

**Definition 3.2** *One says that a set $\Omega \subset \mathbb{R}^d$ is* convex *if and only if, for all $x, y \in \Omega$, the closed segment $[x, y]$ also belongs to $\Omega$.*

*A function $F : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is* convex *if, for all $\lambda \in [0, 1]$ and all $x, y \in \mathbb{R}^d$, one has*

$$F((1 - \lambda)x + \lambda y) \leq (1 - \lambda)F(x) + \lambda F(y). \tag{3.2}$$

*If, whenever the lower bound is not infinite, the inequality above is strict for $\lambda \in (0, 1)$, one says that $F$ is strictly convex.*

Note that, with our definition, convex functions can take the value $+\infty$ but not $-\infty$. In order for the upper-bound to make sense when $F$ takes infinite values, one makes the following convention: $a + (+\infty) = +\infty$ for any $a \in (-\infty, +\infty]$; $\lambda \cdot (+\infty) = +\infty$ for any $\lambda > 0$; $0 \cdot (+\infty)$ is not defined but $0 \cdot (+\infty) + (+\infty) = +\infty$.

**Definition 3.3** *The domain of F, denoted* $\mathrm{dom}(F)$ *is the set of* $x \in \mathbb{R}^d$ *such that* $F(x) < \infty$. *One says that F is* proper *if* $\mathrm{dom}(F) \neq \emptyset$.

*We will only consider proper convex functions in our discussions, which will simply be referred to as convex functions for brevity.*

**Proposition 3.4** *If F is a convex function, then* $\mathrm{dom}(F)$ *is a convex subset of* $\mathbb{R}^d$. *Conversely, if* $\Omega$ *is a convex set and F satisfies* (3.2) *for all* $x, y \in \Omega$ *(i.e., F is convex on* $\Omega$*), then the extension* $\hat{F}$ *defined by* $\hat{F}(x) = F(x)$ *if* $x \in \Omega$ *and* $\hat{F}(x) = +\infty$ *is a convex function defined on* $\mathbb{R}^d$ *(such that* $\mathrm{dom}(\hat{F}) = \Omega$*).*

PROOF The first statement is a direct consequence of (3.2), which implies that $F$ is finite on $[x, y]$ as soon as it is finite at $x$ and at $y$. For the second statement, (3.2) for $\hat{F}$ is true for $x, y \in \Omega$, since it is true for $F$, and the uper-bound is $+\infty$ otherwise. ∎

This proposition shows that there was no real loss of generality in requiring convex functions to be defined on the full space $\mathbb{R}^d$. Note also that the upper bound in (3.2) is infinite unless both $x$ and $y$ belong to $\mathrm{dom}(F)$, so that the inequality only needs to be checked in that case.

One says that a function $F$ is *concave* if and only if $-F$ is convex. All definitions and properties made for convex functions then easily transcribe into similar statements for concave functions. We say that a function $f : I \to (-\infty, +\infty]$ (where $I$ is an interval) is non-decreasing if, for all $x, y \in I$, $x < y$ implies $f(x) \leq f(y)$. We say that $f$ is increasing if if, for all $x, y \in I$, $x < y$ implies $f(x) < f(y)$ if $f(x) < \infty$ and $f(y) = \infty$ otherwise.

Inequality (3.2) has important consequences on minimization problems. For example, it implies the following proposition.

**Proposition 3.5** *Let F be a convex (resp. strictly convex) function on* $\mathbb{R}^d$. *If* $x \in \mathrm{dom}(F)$ *and* $y \in \mathbb{R}^d$, *the function*

$$\lambda \in (0, 1] \mapsto \frac{1}{\lambda}(F((1 - \lambda)x + \lambda y) - F(x)) \tag{3.3}$$

*is non-decreasing (resp. increasing).*

*Conversely, let* $\Omega \subset \mathbb{R}^d$ *be a convex set and* $F : \Omega \to (-\infty, +\infty)$ *be a function such that the expression in* (3.3) *is non-decreasing (resp. increasing) for all* $x \in \mathrm{dom}(F)$ *and* $y \in \mathbb{R}^d$. *Then, the extension* $\hat{F}$ *of F defined in proposition 3.4 is convex (resp. strictly convex).*

PROOF Let $f(\lambda) = (F((1 - \lambda)x + \lambda y) - F(x))/\lambda$. Let $\mu \leq \lambda$ denote $z_\lambda = (1 - \lambda)x + \lambda y$, $z_\mu = (1 - \mu)x + \mu y$. One has $z_\mu = (1 - \nu)x + \nu z_\lambda$, with $\nu = \mu/\lambda$, so that

$$F(z_\mu) \leq (1 - \mu/\lambda)F(x) + (\mu/\lambda)F(z_\lambda).$$

Subtracting $F(x)$ to both sides (which is allowed since $F(x) < \infty$) and dividing by $\mu$ yields

$$f(\mu) \leq f(\lambda).$$

If $F$ is strictly convex, then, either $F(z_\mu) = \infty$, in which case $f(\mu) = f(\lambda) = \infty$, or

$$F(z_\mu) < (1 - \mu/\lambda)F(x) + (\mu/\lambda)F(z_\lambda).$$

as soon as $0 < \mu < \lambda$, yielding

$$f(\mu) < f(\lambda).$$

Now consider the converse statement. By comparing the expression in (3.3) to that obtained with $\lambda = 1$, we find, for all $x, y \in \Omega$

$$\frac{1}{\lambda}(F((1 - \lambda)x + \lambda y) - F(x)) \leq F(y) - f(x)$$

which is (3.2). Since $\hat{F}$ satisfies (3.2) in its domain, it is convex. If the function in (3.3) is increasing, then the inequality is strict for $0 < \lambda < 1$ as soon as the lower bound is finite, and $F$ is strictly convex. ∎

**Corollary 3.6** *If F is convex, any local minimum of F is a global minimum.*

PROOF If $x$ is a local minimum of $F$, then, obviously, $x \in \text{dom}(F)$, and for any $y \in \mathbb{R}^d$ and small enough $\mu > 0$, $F(x) \leq F((1 - \mu)x + \mu y)$. Using the function in (3.3) for $\lambda = \mu$ and for $\lambda = 1$, we get

$$0 \leq \frac{1}{\mu}(F((1 - \mu)x + \mu y) - F(x)) \leq F(y) - F(x)$$

so that $x$ is a global minimum. ∎

### 3.2.3   Relative interior

If $\Omega$ is convex, then $\mathring{\Omega}$ and $\bar{\Omega}$ (its topological interior and closure) are convex too (the easy proof is left to the reader). However, topological interiors of interesting convex sets are often empty, and a more adapted notion of *relative interior* is preferable.

Define the *affine hull* of a set $\Omega$, denoted $\text{aff}(\Omega)$, as the smallest affine subset of $\mathbb{R}^d$ that contains $\Omega$. The vector space parallel to $\text{aff}(\Omega)$ (generated by all differences $x - y$, $x, y \in \Omega$) will be denoted $\overrightarrow{\text{aff}}(\Omega)$. Their dimension $k$, is the largest integer such that there exist $x_0, x_1, \ldots, x_k \in \Omega$ such that $x_1 - x_0, \ldots, x_k - x_0$ are linearly independent. Moreover, given these points, elements of the affine hull are defined through *barycentric coordinates*, yielding

$$\text{aff}(\Omega) = \{x = \lambda^{(0)}x_0 + \cdots + \lambda^{(k)}x_k :, \lambda^{(0)} + \cdots + \lambda^{(k)} = 1\}.$$

The coordinates $(\lambda^{(0)}, \ldots, \lambda^{(k)})$ are uniquely associated to $x \in \text{aff}(\Omega)$ and depend continuously on $x$. They are indeed obtained by solving the linear system

$$x - x_0 = \lambda^{(1)}(x_1 - x_0) + \cdots + \lambda^{(k)}(x_k - x_0)$$

which has a unique solution for $x \in \text{aff}(\Omega)$ by linear independence. To see continuity, one can introduce the $k \times k$ matrix $G$ with entries $G^{(ij)}$ given by the inner products $(x_i - x_0)^T(x_j - x_0)$ and the vector $h(x) \in \mathbb{R}^k$ with entries $h^{(j)}(x) = (x - x_0)^T(x_j - x_0)$. Continuity is then clear since $\lambda = G^{-1}h(x)$.

**Definition 3.7** *If $\Omega$ is a convex set, then its relative interior, denoted* relint$(\Omega)$, *is the set of all $x \in \Omega$ such that there exists $\epsilon > 0$ such that* aff$(\Omega) \cap B(x, \epsilon) \subset \Omega$.

We have the following important property.

**Proposition 3.8** *Let $\Omega$ be a nonempty convex set. If $x \in$ relint$(\Omega)$ and $y \in \Omega$, then $x_\lambda = (1 - \lambda)x + \lambda y \in$ relint$(\Omega)$ for all $\lambda \in [0, 1)$.*

*Moreover* relint$(\Omega)$ *is a nonempty convex set.*

PROOF Take $\epsilon$ such that $B(x, \epsilon) \cap \text{aff}(\Omega) \subset \Omega$. Take any $z \in B(x_\lambda, (1 - \lambda)\epsilon) \cap \text{aff}(\Omega)$. Define $\tilde{z}$ such that $z = (1 - \lambda)\tilde{z} + \lambda y$, i.e.

$$\tilde{z} = \frac{z - \lambda y}{1 - \lambda}.$$

Then $\tilde{z} \in \text{aff}(\Omega)$ and

$$|\tilde{z} - x| = \frac{|z - x_\lambda|}{1 - \lambda} < \epsilon$$

so that $\tilde{z}$, and therefore $z$ belongs to $\Omega$. This proves that $B(x_\lambda, (1 - \lambda)\epsilon) \cap \text{aff}(\Omega) \subset \Omega$ so that $x_\lambda \in$ relint$(\Omega)$.

If both $x$ and $y$ belong to relint$(\Omega)$, then $x_\lambda \in$ relint$(\Omega)$ for $\lambda \in [0, 1]$, showing that this set is convex.

We now show that relint$(\Omega) \neq \emptyset$. Let $k$ be the dimension of aff$(\Omega)$, so that there exist $x_0, x_1, \ldots, x_k \in \Omega$ such that $x_1 - x_0, \ldots, x_k - x_0$ are linearly independent. Consider the "simplex"

$$S = \{\lambda^{(0)}x_0 + \cdots + \lambda^{(k)}x_k :, \lambda^{(0)} + \cdots + \lambda^{(k)} = 1, \lambda^{(j)} \geq 0, j = 0, \ldots, k\},$$

which is included in $\Omega$. Then the average $x = (x_0 + \cdots + x_k)/(k + 1)$ is such that $B(x, \epsilon) \cap \text{aff}(\Omega) \subset S$ for small enough $\epsilon$. Otherwise, there would exist a sequence $y(n) = \lambda^{(0)}(n)x_0 + \cdots + \lambda^{(k)}(n)x_k$ such that $\lambda^{(0)}(n) + \cdots + \lambda^{(k)}(n) = 1$ and at least one $\lambda^{(j)}(n) < 0$ that converges to $x$. Let $y_j$ be the set of elements in this sequence such

that $\lambda^{(j)}(n) < 0$. This set is infinite for at least one $j$ and provides a subsequence of $y$ that also converges to $x$. But this would imply that the $j^{\text{th}}$ barycentric coordinate, which depends continuously on $x$, is non-positive, which is a contradiction.

We therefore have $x \in \text{relint}(\Omega)$, which completes the proof.                    ■

The following proposition provides an equivalent definition of the relative interior.

**Proposition 3.9** *If $\Omega$ is a convex set, then*

$$\text{relint}(\Omega) = \{x \in \Omega : \forall y \in \Omega, \exists \epsilon > 0 \text{ such that } x - \epsilon(y - x) \in \Omega\}. \tag{3.4}$$

So $x$ belongs in the relative interior of $\Omega$ if, for all $y \in \Omega$, the segment $[x, y]$ can be extended on the $x$ side and still remain included in $\Omega$.

Proof  Let $A$ be the set in the r.h.s. of (3.4). The proof that $\text{relint}(\Omega) \subset A$ is straightforward and left to the reader. We consider the reverse inclusion.

Let $x \in A$, and let $y \in \text{relint}(\Omega)$, which is not empty. Then, for some $\epsilon > 0$, we have

$$z = x - \epsilon(y - x) \in \Omega.$$

Since

$$x = \frac{1}{1 + \epsilon}(\epsilon y + z),$$

proposition 3.8 implies that $x \in \text{relint}(\Omega)$.                    ■

Convex functions have important regularity properties in the relative interior of their domain, that we will denote $\text{ridom}(F)$. Importantly:

$$\text{ridom}(F) = \text{relint}(\text{dom}(F)) \neq \text{int}(\text{dom}(F)).$$

A first such property is provided by the next proposition.

**Proposition 3.10** *Let $F$ be a convex function. Then $F$ is locally Lipschitz continuous on* $\text{ridom}(F)$, *i.e., for every compact subset $C \subset \text{ridom}(F)$, there exists a constant $L > 0$ such that $|F(x) - F(y)| \leq L|x - y|$ for all $x, y \in C$.*

This implies, in particular, that $F$ is continuous on $\text{ridom}(F)$.

Proof  Take $x \in \text{ridom}(F)$. Let $K = \left\{h \in \overrightarrow{\text{aff}}(\text{dom}(F)), |h| = 1\right\}$. Then, the segment $[x - ah, x + ah]$ is included in $\text{ridom}(F)$ for small enough $a$ and all $h \in K$. Since $F$ is convex, we have, for $t \leq a$,

$$F(x + th) - F(x) \leq \frac{t}{a}(F(x + ah) - F(x))$$

Writing $x = \lambda(x - ah) + (1 - \lambda)(x + th)$ with $\lambda = t/(t + a)$, we also have

$$F(x) \leq \frac{t}{t + a}(F(x - ah) + \frac{a}{t + a}F(x + th))$$

which can be rewritten as

$$F(x) - F(x + th) \leq \frac{t}{a}(F(x - ah) - F(x)).$$

These two inequalities show that $F$ is continuous at $x$ along any direction in $\overrightarrow{\text{aff}}(\text{dom}(F))$, which implies that $F$ is continuous at $x$. Given this, the differences $F(x + ah) - F(x)$ are bounded over the compact set $C$, by some constant $M$ and, the previous inequalities show that

$$|F(y) - F(x)| \leq \frac{M}{a}|x - y|$$

if $y \in \text{ridom}(F)$, $|y - x| \leq a$. ∎

### 3.2.4   Derivatives of convex functions and optimality conditions

The following theorem provides a stronger version of optimality conditions for the minimization of differentiable convex functions. Note that we have only defined differentiability of functions defined over open sets.

**Theorem 3.11** *Let $F$ be a convex function, with $\text{int}(\text{dom}(F)) \neq \emptyset$. Assume that $x \in \text{int}(\text{dom}(F))$ and that $F$ is differentiable at $x$. Then, for all $y \in \mathbb{R}^d$:*

$$\nabla F(x)^T(y - x) \leq F(y) - F(x). \tag{3.5}$$

*If $F$ is strictly convex, the inequality is strict for $y \neq x$. In particular, $\nabla F(x) = 0$ implies that $x$ is a* global minimizer *of $F$. It is the unique minimizer if $F$ is strictly convex.*

*Conversely, if $F$ is $C^1$ on an open convex set $\Omega$ and satisfies (3.5) for all $x, y \in \Omega$, then F is convex.*

Proof Equation (3.3) implies

$$\frac{1}{\lambda}(F((1 - \lambda)x + \lambda y) - F(x)) \leq F(y) - F(x), 0 < \lambda \leq 1.$$

Taking the limit of the lower bound for $\lambda \to 0$, $\lambda > 0$ yields (3.5). If $F$ is strictly convex, the inequality is strict for $\lambda < 1$ and, since the l.h.s. is increasing in $\lambda$, it remains strict when $\lambda \downarrow 0$.

Conversely, assuming (3.5) for all $x, y \in \Omega$, the derivative of $\lambda \mapsto \frac{1}{\lambda}(F((1 - \lambda)x + \lambda y) - F(x))$ is

$$\frac{1}{\lambda^2}(\lambda \nabla F(x + \lambda h)^T h - F(x + \lambda h) + F(x))$$

with $h = y - x$, which is non-negative by (3.5). This proves that $F$ is convex. If (3.5) holds with a strict inequality, then the derivative is positive and $\frac{1}{\lambda}(F((1 - \lambda)x + \lambda y) - F(x))$ is increasing. ∎

The next proposition describes $C^2$ convex functions in terms of their second derivatives.

**Proposition 3.12** *Let $F$ be convex and twice differentiable at $x \in \text{int}(\text{dom}(F))$. Then $\nabla^2 F(x)$ is positive semi-definite.*

*Conversely, assume that $\Omega = \text{dom}(F)$ is an open set and that $F$ is $C^2$ on $\Omega$ with a positive semi-definite second derivative. Then $F$ (or, rather, its extension $\hat{F}$) is convex. If the second derivative is everywhere positive definite, then $F$ is strictly convex.*

PROOF  Using Taylor formula (1.10) at order 2, we get, for any $h \in \mathbb{R}^d$ with $|h| = 1$,

$$\frac{1}{2}d^2F(x)(h,h) = \frac{1}{2t^2}d^2F(x)(th,th) = \frac{1}{t^2}(F(x+th) - F(x) - t\nabla F(x)^T h) + \epsilon(t) \geq \epsilon(t)$$

with $\epsilon(t) \to 0$ when $t \to 0$, the last inequality deriving from (3.5). This shows that $d^2F(x)(h,h) \geq 0$.

To prove the second statement, assume that $F$ is $C^2$ and $\nabla^2 F$ is positive semi-definite everywhere. Then (1.8) implies

$$F(y) - F(x) - \nabla F(x)^T(y-x) = \frac{1}{2}(y-x)^T \nabla^2 F(z)(y-x)$$

for some $z \in [x,y]$. Since the r.h.s. is non-negative, (3.5) holds. If $\nabla^2 F$ is positive definite everywhere, then the r.h.s. is positive if $y \neq x$ and (3.5) holds with a strict inequality.  ∎

If $F$ is $C^2$ and $\nabla^2 F$ is positive definite and strictly convex, then (1.8) implies that, for some $z \in [x,y]$,

$$F(y) - F(x) - \nabla F(x)^T(y-x) = \frac{1}{2}(y-x)^T \nabla^2 F(z)(y-x) \geq \frac{\rho_{\min}(\nabla^2 F(z))}{2}|y-x|^2$$

where $\rho_{\min}(A)$ denotes the smallest eigenvalue of $A$. If this smallest eigenvalue is bounded from below away from zero, there exists a constant $m > 0$ such that

$$F(y) - F(x) - \nabla F(x)^T(y-x) - \frac{m}{2}|y-x|^2 \geq 0. \tag{3.6}$$

This property is captured by the following definition, which does not require $F$ to be $C^2$.

**Definition 3.13** *A $C^1$ function $F$ is strongly convex if*

1. $\text{int}(\text{dom}(F)) \neq \emptyset$

2. *There exists $m > 0$ such that (3.6) holds for all $x \in \text{int}(\text{dom}(F))$ and $y \in \mathbb{R}^d$.*

We have the following proposition.

**Proposition 3.14** *If F is strongly convex, then it is strictly convex, so that, in particular* argmin *F has at most one element.*

*If* $\text{dom}(F) = \mathbb{R}^d$, *then* argmin *F is not empty.*

PROOF The first part is a direct consequence of (3.6) and theorem 3.11.

For the second part, (3.6) implies that

$$F(x) - F(0) \geq \nabla F(0)^T x + \frac{m}{2}|x|^2 \geq |x|\left(\frac{m}{2}|x|^2 - |\nabla F(0)|\right)$$

This shows that $F(x) > F(0)$ if $|x| > 2|\nabla F(0)|/m := r$ so that

$$\text{argmin}\, F = \underset{\bar{B}(0,r)}{\text{argmin}}\, F.$$

The set in the r.h.s. involves the minimization of a continuous function on a compact set, and is therefore not empty. ■

We will use the following definition.

**Definition 3.15** *A function* $F : \Omega \to \mathbb{R}^m$ *is* $L$-$C^k$, $L$ *being a positive number, if it is* $C^k$ *and*
$$|d^k F(x) - d^k F(y)| \leq L|x - y|.$$

If $F$ is $L$-$C^k$, then Taylor formula ((1.9)) implies

$$\left| f(x+h) - f(x) - df(x)h - \frac{1}{2}d^2 f(x)(h,h) - \cdots - \frac{1}{k!}d^k f(x)(h,\ldots,h) \right| \leq \frac{L|h|^{k+1}}{(k+1)!} \quad (3.7)$$

for which we used the fact that

$$\int_0^1 t(1-t)^{k-1}dt = \int_0^1 (1-t)^{k-1}dt - \int_0^1 (1-t)^k dt = \frac{1}{k} - \frac{1}{k+1} = \frac{1}{k(k+1)}.$$

If $F$ is strongly convex and is, in addition, $L$-$C^1$ for some $L$, then using (3.7), one gets the double inequality, for all $x, y \in \text{int}(\text{dom}(F))$:

$$\frac{m}{2}|y - x|^2 \leq F(y) - F(x) - \nabla F(x)^T(y - x) \leq \frac{L}{2}|y - x|^2. \quad (3.8)$$

The following proposition will be used later.

**Proposition 3.16** *Assume that F is strongly convex, satisfying* (3.6), *and that* $\arg\min F =$ $\{x^*\}$ *with* $x^* \in \text{int}(\text{dom}(F))$. *Then, for all* $x \in \text{int}(\text{dom}(F))$:

$$\frac{m}{2}|x - x^*|^2 \leq F(x) - F(x^*) \leq \frac{1}{2m}|\nabla F(x)|^2 \tag{3.9}$$

Proof  Since $\nabla F(x^*) = 0$, the first inequality is a consequence of (3.6) applied to $x = x^*$. Switching the role of $x$ and $x^*$, we have

$$F(x^*) - F(x) - \nabla F(x)^T(x^* - x) \geq \frac{m}{2}|x - x^*|^2$$

so that

$$0 \leq F(x) - F(x^*) \leq -\nabla F(x)^T(x^* - x) - \frac{m}{2}|x - x^*|^2 \leq |\nabla F(x)||x - x^*| - \frac{m}{2}|x - x^*|^2 \tag{3.10}$$

The maximum of the r.h.s. with respect to $|x - x^*|$ is attained at $|\nabla F(x)|/m$, showing that

$$F(x) - F(x^*) \leq \frac{1}{2m}|\nabla F(x)|^2,$$

which is the second inequality.                                                        ∎

### 3.2.5   Direction of descent and steepest descent

Gradient-based algorithms for optimization iteratively update the variable $x$, creating a sequence governed by an equation taking the form $x_{t+1} = x_t + \alpha_t h_t$ with $\alpha_t > 0$ and $h_t \in \mathbb{R}^d$. To ensure that the objective function $F$ decreases at each step, $h_t$ is chosen to be a direction of descent for $F$ at $x_t$, a notion which, as seen below, is closely connected with the direction of $\nabla F(x_t)$.

**Definition 3.17** *Let* $\Omega$ *be open in* $\mathbb{R}^d$ *and* $F : \Omega \to \mathbb{R}$ *be a* $C^1$ *function. A direction of descent for F at* $x \in \Omega$ *is a vector* $h \neq 0 \in \mathbb{R}^d$ *such that there exists* $\epsilon_0 > 0$ *such that* $F(x + \epsilon h) < F(x)$ *for all* $\epsilon \in (0, \epsilon_0]$.

**Proposition 3.18** *Assume that* $F : \Omega \to \mathbb{R}$ *is* $C^1$ *and take* $x \in \Omega$. *Then any direction h such that* $h^T \nabla F(x) < 0$ *is a direction of descent for F at x. Conversely, if h is a direction of descent, then* $h^T \nabla F(x) \leq 0$.

Proof  We have the first-order expansion $F(x+\epsilon h)-F(x) = \epsilon h^T \nabla F(x)+o(\epsilon)$. If $h^T \nabla F(x) < 0$, the r.h.s. is negative for small enough $\epsilon$ and $h$ is a direction of descent. Similarly, if $h^T \nabla F(x) > 0$, the r.h.s. is positive for small enough $\epsilon$ and $h$ cannot be a direction of descent.                                                        ∎

In particular, $h = -\nabla F(x)$ is always a direction of descent. It is called the steepest descent direction because it minimizes $h \mapsto \partial_\alpha F(x + \alpha h)|_{\alpha=0}$ over all $h$ such that $|h|^2 = 1$. However, this designation has a character of optimality that may be misleading, because using the Euclidean norm for the condition $|h|^2 = 1$ is not necessarily adapted to the optimization problem at hand. In the absence of additional information on the problem, it does have a canonical nature, as it is (up to rescaling) the only norm invariant to rotations (including permutations) of the coordinates. Such invariance is not necessarily desirable when the variable $x$ has a known structure (e.g., it is organized on a graph) which would be broken by permutation. Also, steepest refers to a local "greedy" evaluation, but may not be optimal from a global perspective. A simple example to illustrate this is the case of a quadratic function

$$F(x) = \frac{1}{2}x^T A x - b^T x$$

where $A \in \mathcal{S}_n^{++}$ is a positive definite symmetric matrix. Then $\nabla F(x) = Ax - b$, but one may argue that $\nabla_A F(x) = A^{-1}\nabla F(x)$ (defined in (1.3)) is a better choice, because it allows the algorithm to reach the minimizer of $F$ in one step, since $x - \nabla_A F(x) = A^{-1}b$ (this statement disregards the cost associated in solving the system $Ax = b$, which can be an important factor in large dimension). Importantly, if $F$ is any $C^1$ function, and $A \in \mathcal{S}_n^{++}$, the minimizer of $h \mapsto \partial_\alpha F(x + \alpha h)|_{\alpha=0}$ over all $h$ such that $h^T A h = 1$ is given by $-\nabla_A F(x)$, i.e., $-\nabla_A F(x)$ is the steepest descent for the norm associated with $A$. This yields a general version of steepest descent methods, iterating

$$x_{t+1} = x_t - \alpha_t \nabla_{A_t} F(x_t)$$

with $\alpha_t > 0$ and $A_t \in \mathcal{S}_n^{++}$.

One can also notice that $\nabla_A F(x)$ is also a minimizer of

$$F(x) + \nabla F(x)^T h + \frac{1}{2}h^T A h.$$

When $\nabla^2 F(x)$ is positive definite, it is then natural to choose it as the matrix $A$, therefore taking $h = -\nabla^2 F(x)^{-1}\nabla F(x)$. This provides Newton's method for optimization. However, Newton method requires computing second derivatives of $F$, which can be computationally costly. It is, moreover, not a gradient-based method, which is the focus of this discussion.

### 3.2.6 Convergence

We now consider a descent algorithm

$$x_{t+1} = x_t + \alpha_t h_t \tag{3.11}$$

where $h_t$ is a direction of descent at $x_t$ for the objective function $F$. To ensure convergence, suitable choices for the direction of descent and the step must be made at each iteration, and some assumptions on the objective function are needed.

Regarding the direction of descent, which must satisfy $h_k^T \nabla F(x_k) \leq 0$, we will assume a uniform control away from orthogonality to the gradient, with the condition

$$-h_t^T \nabla F(x_t) \geq \epsilon |h_t| |\nabla F(x_t)| \tag{3.12a}$$

for some fixed $\epsilon > 0$. Without loss of generality (given that a multiplicative step $\alpha_t$ must also be chosen), we assume that $h_t$ is commensurable to the gradient, namely, that

$$\gamma_1 |\nabla F(x_t)| \leq |h_t| \leq \gamma_2 |\nabla F(x_t)| \tag{3.12b}$$

for fixed $0 < \gamma_1 \leq \gamma_2$. If $h_t = \nabla_{A_t} F$, these assumptions are satisfied as soon as the smallest and largest eigenvalues of $A_t$ are controlled along the trajectory.

We have the following proposition.

**Proposition 3.19** *Assume that $F$ is $L$-$C^1$. Assume that $x_t$ satisfies (3.11) and that (3.12a) and (3.12b) hold. Then, there exist constants $\bar{\alpha} > 0$ and $C > 0$ that depends on $\gamma_1, \gamma_2$ and $\epsilon$, such that, for $\alpha_t \leq \bar{\alpha}$, one has*

$$F(x_{t+1}) - F(x_t) \leq -C\alpha_t |\nabla F(x_t)|^2. \tag{3.13}$$

Proof Applying (3.7) to $x_t$ and $x_{t+1}$, we get

$$F(x_{t+1}) - F(x_t) - \alpha_t \nabla F(x_t)^T h_t \leq \frac{L}{2} \alpha_t^2 |h_t|^2$$

Using (3.11) and (3.12a), this gives

$$F(x_{t+1}) - F(x_t) + \alpha_t \epsilon \gamma_1 |\nabla F(x_t)|^2 \leq \frac{L}{2} \alpha_t^2 \gamma_2^2 |\nabla F(x_t)|^2$$

so that

$$F(x_{t+1}) - F(x_t) \leq -\alpha_t \left( \epsilon \gamma_1 - \alpha_t \gamma_2^2 L/2 \right) |\nabla F(x_t)|^2.$$

It suffices to take $\bar{\alpha} = \epsilon \gamma_1 / L \gamma_2^2$ and $C = \epsilon \gamma_1 / 2$ to obtain (3.13). ∎

Iterating (3.13) for $t = 1$ to $t = T - 1$ yields

$$\sum_{t=1}^{T} \alpha_t |\nabla F(x_t)|^2 \leq \frac{1}{C}(F(x_1) - F(x_T)).$$

If $F$ is bounded from below, and one takes $\alpha_t = \bar{\alpha}$ for all $t$, one deduces that

$$\min \left\{ |\nabla F(x_t)|^2 : t = 1, \ldots, T \right\} \leq \frac{F(x_1) - \inf F}{CT\bar{\alpha}}.$$

We can deduce from this, for example, that there exists a sequence $t_1 < \cdots < t_n < \cdots$ such that $\nabla F(x_{t_k}) \to 0$ when $k \to \infty$. In particular, if one runs (3.11) until $|\nabla F(x_t)|$ is smaller than a given tolerance level (which is standard), the procedure is guaranteed to terminate in a finite number of steps.

Stronger results may be obtained under stronger assumptions on $F$ and on the algorithm. The first assumption is an inequality similar to (3.13) and requires that, for some constant $C > 0$,

$$F(x_{t+1}) - F(x_t) \leq -C|\nabla F(x_t)|^2. \tag{3.14}$$

Such an inequality can be deduced from (3.13) under the additional assumption that $\alpha_t$ is bounded from below and we will discuss later line search strategies that ensure its validity. The second assumption is that $F$ is convex.

**Theorem 3.20** *Assume that $F$ is convex and finite and that its sub-level set $[F \leq F(x_0)]$ is bounded. Assume that $\arg\min F$ is not empty and let $x^*$ be a minimizer of $F$. If (3.14) is true, then*

$$F(x_t) - F(x_*) \leq \frac{R^2}{C(t+1)}$$

*with $R = \max\{|x - x_*| : F(x) \leq F(x_0)\}$.*

Proof Note that the algorithm never leaves $[F \leq F(x_0)]$. We have

$$F(x_{t+1}) - F(x^*) \leq F(x_t) - F(x^*) - C|\nabla F(x_t)|^2.$$

Moreover, by convexity, $F(x^*) - F(x_t) \geq \nabla F(x_t)^T (x^* - x_t)$, so that

$$F(x_t) - F(x^*) \leq \nabla F(x_t)^T (x_t - x^*) \leq |\nabla F(x_t)|R.$$

Combining these two inequalities, we get

$$F(x_{t+1}) - F(x^*) \leq F(x_t) - F(x^*) - \frac{C}{R^2}(F(x_t) - F(x^*))^2.$$

Introducing $\delta_t = (C/R^2)(F(x_t) - F(x^*))$, this inequality implies

$$\delta_{t+1} \leq \delta_t(1 - \delta_t).$$

Taking inverses, we get

$$\frac{1}{\delta_{t+1}} \geq \frac{1}{\delta_t} + \frac{1}{1 - \delta_t} \geq \frac{1}{\delta_t} + 1$$

which implies $\frac{1}{\delta_t} \geq t + 1$ or $\delta_t \leq 1/(t+1)$), which in turn implies the statement of the theorem. ∎

A faster convergence rate can be obtained if $F$ is assumed to be strongly convex. Indeed, if (3.6) and (3.14) are satisfied, then (using proposition 3.16),

$$
\begin{aligned}
F(x_{t+1}) - F(x^*) &\leq F(x_t) - F(x^*) - C|\nabla F(x_t)|^2 \\
&\leq F(x_t) - F(x^*) - 2Cm(F(x_t) - F(x^*)) \\
&= (1 - 2Cm)(F(x_t) - F(x^*)).
\end{aligned}
$$

We therefore get the proposition:

**Proposition 3.21** *If $F$ is finite and satisfies* (3.6), *and if the descent algorithm satisfies* (3.14), *then*
$$
F(x_t) - F(x^*) \leq (1 - 2Cm)^t (F(x_0) - F(x^*)).
$$

### 3.2.7  Line search

Proposition 3.19 states that, to ensure that (3.14) holds, it suffices to take a small enough step parameter $\alpha$. However, the values of $\alpha$ that are acceptable depend on properties of the objective function that are rarely known in practice. Moreover, even if a valid choice is determined (this can sometimes be done in practice by trial and error), setting a fixed value of $\alpha$ for the whole algorithm is often too conservative, as the best $\alpha$ when starting the algorithm may be different from the best one close to convergence.

For this reason, most gradient descent procedures select a parameter $\alpha_t$ at each step using a line search. Given a current position and direction of descent $h$, a line search explores the values of $F(x + \alpha h)$, $\alpha \in (0, \alpha_{\max}]$ in order to discover some $\alpha^*$ that satisfies some desirable properties. We will assume in the following that $x$ and $h$ satisfy (3.12a) and (3.12b) for fixed $\epsilon, \gamma_1, \gamma_2$.

One possible strategy is to define $\alpha^*$ as a minimizer of the scalar function

$$
f_h(\alpha) = F(x + \alpha h)
$$

over $(0, \alpha_{\max}]$ for a given upper-bound $\epsilon_{\max}$. This can be implemented using, e.g., binary or ternary search algorithms, but such algorithms would typically require a large number of number of evaluations of the function $F$, and would be too costly to be run at each iteration of a gradient descent procedure.

Based on the previous convergence study, we should be happy with a line search procedure that ensures that (3.14) is satisfied for some fixed value of the constant $C$. One such condition is the so-called Armijo rule that requires (with a fixed, typically small, value of $c_1 > 0$):

$$
f_h(\alpha) \leq f_h(0) + c_1 \alpha h^T \nabla f(x). \tag{3.15}
$$

We know that, under the assumptions of proposition 3.19, this condition can always be satisfied with a small enough value of $\alpha$. Such a value can be determined using a "backtracking procedure," which, given $\alpha_{\max}$ and $\rho \in (0,1)$, takes $\alpha = \rho^k \alpha_{\max}$ where $k$ is the smallest integer such that (3.15) is satisfied. This value of $k$ is then determined iteratively, trying $\alpha_{\max}$, $\rho\alpha_{\max}$, $\rho^2\alpha_{\max}, \dots$ until (3.15) is true (this provides the "backtracking method").

A stronger requirement in the line search is to ensure that $\partial f_h(\alpha)$ is not "too negative" since one would otherwise be able to further reduce $f_h$ by taking a larger value of $\alpha$. This leads to the weak Wolfe conditions, which combine the Armijo's rule in (3.15) and

$$\partial f_h(\alpha) = h^T \nabla F(x + \alpha h) \geq c_2 h^T \nabla F(x) \tag{3.16a}$$

for some constant $c_2 \in (c_1, 1)$. The strong Wolfe conditions require (3.15) and

$$|h^T \nabla F(x + \alpha h)| \leq c_2 |h^T \nabla F(x)|. \tag{3.16b}$$

(Since $h$ is a direction of descent, (3.16b) requires (3.16a) and the fact that $h^T \nabla F(x + \alpha h)$ does not take too large positive values.) If $F$ is $L$-$C^1$, these conditions, with (3.12a) and (3.12b), imply (3.14). Indeed, (3.16a) and the $L$-$C^1$ condition imply

$$-(1 - c_2)h^T \nabla F(x) \leq h^T (\nabla F(x + \alpha h) - \nabla F(x)) \leq L\alpha |h|^2$$

and (3.12a) and (3.12b) give

$$(1 - c_2)\epsilon |\nabla F(x)|^2 \leq \alpha L \gamma_2^2 |\nabla F(x)|^2$$

showing that $\alpha \geq (1 - c_2)\epsilon/(L\gamma_2^2)$. Moreover

$$F(x + \alpha h) \leq F(x) + c_1 \alpha h^T \nabla f(x) \leq F(x) - c_1 \alpha \epsilon |\nabla F(x)|^2$$

so that

$$F(x + \alpha h) \leq F(x) - \frac{c_1(1 - c_2)\epsilon^2}{L\gamma_2^2} |\nabla F(x)|^2.$$

We have just proved the following proposition.

**Proposition 3.22** *Assume that $F$ is $L$-$C^1$ and that (3.12a), (3.12b), (3.15) and (3.16a) are satisfied. Then there exists $C > 0$, depending only of $L, \epsilon, \gamma_2, c_1$ and $c_2$such that*

$$F(x + \alpha h) \leq F(x) - C|\nabla F(x)|^2.$$

The Wolfe conditions can always be satisfied by some $\alpha$ as soon as $F$ is $C^1$ and bounded from below, and $h^T \nabla F(x) < 0$. The next proposition shows this result for the weak condition, while providing an algorithm finding an $\alpha$ that satisfies it in a finite number of steps.

**Proposition 3.23** *Let $f : \alpha \mapsto f(\alpha)$ be a $C^1$ function defined on $[0, +\infty)$ such that $f$ is bounded from below and $\partial_\alpha f(0) < 0$. Let $0 < c_1 < c_2 < 1$.*

*Let $\alpha_{0,0} = \alpha_{0,1} = 0$ and $\alpha_0 > 0$. Define recursively sequences $\alpha_{n,0}, \alpha_{n,1}$ and $\alpha_n$ as follows.*

*(i)  If $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$ and $\partial f(\alpha_n) \geq c_2 \partial_\alpha f(0)$ stop the construction.*

*(ii)  If $f(\alpha_n) > f(0) + c_1 \alpha_n \partial_\alpha f(0)$ let $\alpha_{n+1} = (\alpha_n + \alpha_{n,0})/2$, $\alpha_{n+1,1} = \alpha_n$ and $\alpha_{n+1,0} = \alpha_{n,0}$.*

*(iii)  If $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$ and $\partial f(\alpha_n) < c_2 \partial_\alpha f(0)$:*

    *(a) If $\alpha_{n,1} = 0$, let $\alpha_{n+1} = 2\alpha_n$, $\alpha_{n+1,0} = \alpha_n$ and $\alpha_{n+1,1} = \alpha_{n,1}$.*

    *(b) If $\alpha_{n,1} > 0$, let $\alpha_{n+1} = (\alpha_n + \alpha_{n,1})/2$, $\alpha_{n+1,0} = \alpha_n$ and $\alpha_{n+1,1} = \alpha_{n,1}$.*

*Then the sequences are always finite, i.e., the algorithm terminates in a finite number of steps.*

PROOF  Assume, to get a contradiction, that the algorithm runs indefinitely, so that case (i) never occurs. If case (ii) never occurs, then one runs step (iii-a) indefinitely, so that $\alpha_n \to \infty$ with $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$, and $f$ cannot be bounded from below, yielding a contradiction. As soon as case (ii) occurs, we have, at every step, $\alpha_{n,0} \geq \alpha_{n-1,0}$, $\alpha_{n,1} \leq \alpha_{n-1,1}$, $\alpha_n \in [\alpha_{n,0}, \alpha_{n,1}]$, $f(\alpha_{n,1}) > f(0) + c_1 \alpha_{n,1} \partial_\alpha f(0)$, $f(\alpha_{n,0}) \leq f(0) + c_1 \alpha_{n,0} \partial_\alpha f(0)$ and $\partial f(\alpha_{n,0}) < c_2 \partial_\alpha f(0)$. This implies that

$$f(\alpha_{n,1}) - f(\alpha_{n,0}) > c_1 (\alpha_{n,1} - \alpha_{n,0}) \partial_\alpha f(0).$$

Moreover, the updates imply that $(\alpha_{n+1,1} - \alpha_{n+1,0}) = (\alpha_{n,1} - \alpha_{n,0})/2$. This requires that the three sequences $\alpha_n, \alpha_{n,0}$ and $\alpha_{n,1}$ converge to the same limit, $\alpha$. We have

$$\partial_\alpha f(\alpha) = \lim_{n \to \infty} \frac{f(\alpha_{n,1}) - f(\alpha_{n,0})}{\alpha_{n,1} - \alpha_{n,0}} \geq c_1 \partial_\alpha f(0)$$

and

$$\partial_\alpha f(\alpha) = \lim_{n \to \infty} \partial_\alpha f(\alpha_{n,0}) \leq c_2 \partial_\alpha f(0)$$

yielding $c_1 \partial_\alpha f(0) \leq c_2 \partial_\alpha f(0)$ which is impossible since $c_2 > c_1$ and $\partial_\alpha f(0) < 0$.  ∎

The existence of $\alpha$ satisfying the strong Wolfe condition is a consequence of the following proposition, which also provides an algorithm.

**Proposition 3.24** *Let $f : \alpha \mapsto f(\alpha)$ be a $C^1$ function defined on $[0, +\infty)$ such that $f$ is bounded from below and $\partial_\alpha f(0) < 0$. Let $0 < c_1 < c_2 < 1$.*

*Let $\alpha_{0,0} = \alpha_{0,1} = 0$ and $\alpha_0 > 0$. Define recursively sequences $\alpha_{n,0}, \alpha_{n,1}$ and $\alpha_n$ as follows.*

*(i) If $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$ and $|\partial_\alpha f(\alpha_n)| \leq c_2 |\partial_\alpha f(0)|$ stop the construction.*

*(ii) If $f(\alpha_n) > f(0) + c_1 \alpha_n \partial_\alpha f(0)$ let $\alpha_{n+1} = (\alpha_n + \alpha_{n,0})/2$, $\alpha_{n+1,1} = \alpha_n$ and $\alpha_{n+1,0} = \alpha_{n,0}$.*

*(iii) If $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$ and $|\partial_\alpha f(\alpha_n)| > c_2 |\partial_\alpha f(0)|$:*

*(a) If $\alpha_{n,1} = 0$ and $\partial_\alpha f(\alpha_n) > -c_2 \partial_\alpha f(0)$, let $\alpha_{n+1} = 2\alpha_n$, $\alpha_{n+1,0} = \alpha_{n,0}$ and $\alpha_{n+1,1} = \alpha_{n,1}$.*

*(b) If $\alpha_{n,1} = 0$ and $\partial_\alpha f(\alpha_n) < c_2 \partial_\alpha f(0)$, let $\alpha_{n+1} = 2\alpha_n$, $\alpha_{n+1,0} = \alpha_n$ and $\alpha_{n+1,1} = \alpha_{n,1}$.*

*(c) If $\alpha_{n,1} > 0$ and $\partial_\alpha f(\alpha_n) > -c_2 \partial_\alpha f(0)$, let $\alpha_{n+1} = (\alpha_n + \alpha_{n,0})/2$, $\alpha_{n+1,1} = \alpha_n$ and $\alpha_{n+1,0} = \alpha_{n,0}$.*

*(d) If $\alpha_{n,1} > 0$ and $\partial_\alpha f(\alpha_n) < c_2 \partial_\alpha f(0)$, let $\alpha_{n+1} = (\alpha_n + \alpha_{n,1})/2$, $\alpha_{n+1,0} = \alpha_n$ and $\alpha_{n+1,1} = \alpha_{n,1}$.*

*Then the sequences are always finite, i.e., the algorithm terminates in a finite number of steps.*

PROOF Assume that the algorithm runs indefinitely in order to get a contradiction. If the algorithm never enters case (ii), then $\alpha_{n,1} = 0$ for all $n$, $\alpha_n$ tends to infinity and $f(\alpha_n) \leq f(0) + c_1 \alpha_n \partial_\alpha f(0)$, which contradicts the fact that $f$ is bounded from below.

As soon as the algorithm enter (ii), we have, for all subsequent iterations: $\alpha_{n,0} \leq \alpha_n \leq \alpha_{n,1}$, $\alpha_{n+1,0} \geq \alpha_{n,0}$, $\alpha_{n+1,1} \leq \alpha_{n,1}$ and $\alpha_{n+1,1} - \alpha_{n+1,0} = (\alpha_{n,1} - \alpha_{n,0})/2$. This implies that both $\alpha_{n,0}$ and $\alpha_{n,1}$ converge to the same limit $\alpha$.

Moreover, we have, at each step:

$$f(\alpha_{n,1}) > f(0) + c_1 \alpha_{n,1} \partial_\alpha f(0) \text{ or } \partial_\alpha f(\alpha_{n,1}) > -c_2 \partial_\alpha f(0)$$

and

$$f(\alpha_{n,0}) \leq f(0) + c_1 \alpha_{n,0} \partial_\alpha f(0) \text{ and } \partial_\alpha f(\alpha_{n,0}) \leq c_2 \partial_\alpha f(0).$$

This implies that, at each step:

$$\frac{f(\alpha_{n,1}) - f(\alpha_{n,0})}{\alpha_{n,1} - \alpha_{n,0}} > c_1 \partial_\alpha f(0) \text{ or } \partial_\alpha f(\alpha_{n,1}) > -c_2 \partial_\alpha f(0)$$

and

$$\partial_\alpha f(\alpha_{n,0}) \leq c_2 \partial_\alpha f(0).$$

There inequalities remain satisfied at the limit, and we must have

$$\partial_\alpha f(\alpha) > c_1 \partial_\alpha f(0) \text{ or } \partial_\alpha f(\alpha) > -c_2 \partial_\alpha f(0)$$

and

$$\partial_\alpha f(\alpha) \leq c_2 \partial_\alpha f(0),$$

which is a contradiction since $c_2 > c_1$ and $\partial_\alpha f(0) < 0$. ∎

## 3.3   Stochastic gradient descent

### 3.3.1   Stochastic approximation methods

In some situations, the computation of $\nabla F$ can be too costly, if not intractable, to run gradient descent updates while a low-cost stochastic approximation is available. For example, if $F$ is an average of a sum of many terms, the approximation may simply be based on averaging over a randomly selected subset of the terms. This leads to a *stochastic approximation algorithm* [163, 113, 25, 67] called stochastic gradient descent (SGD).

A general stochastic approximation algorithm of the Robbins-Monro type updates a parameter, denoted $x \in \mathbb{R}^d$, using stochastic rules. One associates to each $x$ a probability distribution $(\pi_x)$ on some set $\mathcal{S}$, and, for some function $H : \mathbb{R}^d \times \mathcal{S} \to \mathbb{R}^d$, considers the sequence of random iterations:

$$\begin{cases} \xi_{t+1} \sim \pi_{X_t} \\ X_{t+1} = X_t + \alpha_{t+1} H(X_t, \xi_{t+1}) \end{cases} \tag{3.17}$$

where $\xi_{t+1}$ is a random variable and the notation $\xi_{t+1} \sim \pi_{X_t}$ should be interpreted as the more precise statement that the conditional distribution of $\xi_{t+1}$ given all past random variables $\mathcal{U}_t = (\xi_1, X_1, \dots, \xi_t, X_t)$ only depends on $X_t$ and is given by $\pi_{X_t}$.

It is sometimes assumed in the literature that $\pi_x$ does not depend on $x$. This is no real loss of generality because under mild assumptions, a random variable $\xi$ following $\pi_x$ can be generated as function $U(x, \tilde{\xi})$ where $\tilde{\xi}$ follows a fixed distribution (such as that of a family of independent uniformly distributed variables) and one can replace $H(x, \xi)$ by $H(x, U(x, \tilde{\xi}))$. On the other hand, allowing $\pi$ to depend on $x$ brings little additional complication in the notation, and corresponds to the natural form of many applications.

More complex situations can also be considered, in which $\xi_{t+1}$ is not conditionally independent of the past variables given $X_t$. For example, the conditional distribution of $\xi_{t+1}$ given the past may also depends on $\xi_t$, which allows for the combination of stochastic gradient methods with Markov chain Monte-Carlo methods. This situation is studied, for example, in Métivier and Priouret [138], Benveniste et al. [25], and we will discuss an example in section 17.2.2.

### 3.3.2   Deterministic approximation and convergence study

Introduce the function

$$\bar{H}(x) = E_{\pi_x}(H(x, \cdot))$$

and write

$$X_{t+1} = X_t + \alpha_{t+1}\bar{H}(X_t) + \alpha_{t+1}\eta_{t+1}$$

with $\eta_{t+1} = H(X_t, \xi_{t+1}) - \bar{H}(X_t)$ in order to represent the evolution of $X_t$ in (3.17) as a perturbation of the deterministic algorithm

$$\bar{x}_{t+1} = \bar{x}_t + \alpha_{n+1}\bar{H}(\bar{x}_t) \tag{3.18}$$

by the "noise term" $\alpha_{t+1}\eta_{t+1}$. In many cases, the deterministic algorithm provides the limit behavior of the stochastic sequence, and one should ensure that this limit is as desired. By definition, the conditional expectation of $\eta_{t+1}$ given $\mathcal{U}_t$ (the past) is zero and one says that $\alpha_{t+1}\eta_{t+1}$ is a "martingale increment." Then,

$$M_T = \sum_{t=0}^{T} \alpha_{t+1}\eta_{t+1} \tag{3.19}$$

is called a "martingale." The theory of martingales offers numerous tools for controlling the size of $M_T$ and is often a key element in proving the convergence of the method.

Many convergence results have been provided in the literature and can be found in textbooks or lecture notes such as Benaïm [23], Kushner and Yin [113], Benveniste et al. [25]. These results rely on some smoothness and growth assumptions made on the function $H$, and on the dynamics of the deterministic equation (3.18). Depending on these assumptions, proofs may become quite technical. We will here restrict to a reasonably simple context and assume that

(H1) There exists a constant $C$ such that, for all $x \in \mathbb{R}^d$,

$$E_{\pi_x}(|H(x,\cdot)|^2) \leq C(1 + |x|^2).$$

(H2) There exists $x^* \in \mathbb{R}^d$ and $\mu > 0$ such that, for all $x \in \mathbb{R}^d$

$$(x - x^*)^T \bar{H}(x) \leq -\mu|x - x^*|^2.$$

Assuming this, let $A_t = |X_t - x^*|^2$ and $a_t = \mathbb{E}(A_t)$. Then, using (3.17),

$$A_{t+1} = A_t + 2\alpha_{t+1}(X_t - x^*)^T H(X_t, \xi_{t+1}) + \alpha_{t+1}^2 |H(X_t, \xi_{t+1})|^2.$$

Taking the conditional expectation given past variables yields

$$\begin{aligned}
\mathbb{E}(A_{t+1} \mid \mathcal{U}_t) &= A_t + 2\alpha_{t+1}(X_t - x^*)^T \bar{H}(X_t) + \alpha_{t+1}^2 E_{\pi_{x_t}}(|H(X_t, \cdot)|^2) \\
&\leq A_t - 2\alpha_{t+1}\mu A_t + \alpha_{t+1}^2 C(1 + |X_t|^2) \\
&\leq (1 - 2\alpha_{t+1}\mu + C\alpha_{t+1}^2)A_t + \alpha_{t+1}^2 \tilde{C}
\end{aligned}$$

with $\tilde{C} = 1 + |x^*|^2$. Taking expectations on both sides yields

$$a_{t+1} \le (1 - 2\alpha_{t+1}\mu + C\alpha_{t+1}^2)a_t + \alpha_{t+1}^2 \tilde{C}. \tag{3.20}$$

We state the next step in the computation as a lemma.

**Lemma 3.25** *Assume that the sequence $a_t$ satisfies the recursive inequality*

$$a_{t+1} \le (1 - \delta_t)a_t + \epsilon_t \tag{3.21}$$

*with $0 \le \delta_t \le 1$. Let $v_{k,t} = \prod_{j=k+1}^{t}(1 - \delta_j)$. Then*

$$a_t \le a_0 v_{0,t} + \sum_{k=1}^{t} \epsilon_k v_{k,t}. \tag{3.22}$$

PROOF  Letting $b_t = a_t/v_{0,t}$, we get

$$b_{t+1} \le b_t + \frac{\epsilon_{t+1}}{v_{0,t+1}}$$

so that

$$b_t \le b_0 + \sum_{k=1}^{t} \frac{\epsilon_k}{v_{0,k}},$$

and

$$a_t \le a_0 v_{0,t} + \sum_{k=1}^{t} \epsilon_k v_{k,t}. \qquad \blacksquare$$

Using (3.20), we can apply this lemma with $\epsilon_t = \tilde{C}\alpha_t^2$ and $\delta_t = 2\alpha_t\mu - C\alpha_t^2$, making the additional assumption that, for all $t$, $\alpha_t < \min(\frac{1}{2\mu}, \frac{2\mu}{C})$, which ensures that $0 < \delta_t < 1$.

Starting with a simple case, assume that the steps $\gamma_t$ are constant, equal to some value $\gamma$ (yielding also constant $\delta$ and $\epsilon$). Then, (3.22) gives

$$a_t \le a_0(1 - \delta)^t + \epsilon \sum_{k=1}^{t} (1 - \delta)^{t-k-1} \le a_0(1 - \delta)^t + \frac{\epsilon}{\delta}. \tag{3.23}$$

Returning to the expression of $\delta$ and $\epsilon$ as functions of $\alpha$, this gives

$$a_t \le a_0(1 - 2\alpha\mu + \alpha^2 C)^t + \frac{\alpha\tilde{C}}{2\mu - \alpha C}.$$

This shows that $\limsup a_t = O(\alpha)$.

Return to the general case in which the steps depend on $t$, we will use the following simple result, that we state as a lemma for future reference.

**Lemma 3.26** *Assume that the double indexed sequence $w_{st}$, $s \le t$ of non-negative numbers is bounded and such that, for all $s$, $\lim_{t \to \infty} w_{st} = 0$. Let $\beta_1, \beta_2, \dots$ be such that*

$$\sum_{t=1}^{\infty} |\beta_t| < \infty.$$

*Then*

$$\lim_{t \to \infty} \sum_{s=1}^{t} \beta_s w_{st} = 0.$$

PROOF For any $t_0$, we have

$$\left| \sum_{s=1}^{t} \beta_s w_{st} \right| \le \max_s |\beta_s| \sum_{s=1}^{t_0} w_{st} + \max_{s,t} |w_{st}| \sum_{s=t_0+1}^{\infty} |\beta_s|$$

so that

$$\limsup_{t \to \infty} \left| \sum_{s=1}^{t} \beta_s w_{st} \right| \le \max_{s,t} |w_{st}| \sum_{s=t_0+1}^{\infty} |\beta_s|$$

and since this upper bound can be made arbitrarily small, the result follows. ∎

Lemma 3.25 implies that

$$a_t \le a_0 v_{0,t} + \tilde{C} \sum_{s=1}^{t} \alpha_{s+1}^2 v_{s,t}.$$

Assume that

(H3) $\sum_{k=1}^{\infty} \alpha_k = \infty$ and $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$,

Then $\lim_{t \to \infty} v_{st} = 0$ for all $s$ and lemma 3.26 implies that $a_t$ tends to zero. So, we have just proved that, if (H1), (H2) and (H3) are true, the sequence $X_t$ converges in the $L^2$ sense to $x^*$. Actually, under these conditions, one can show that $X_t$ converge to $x^*$ almost surely, and we refer to Benveniste et al. [25], Chapter 5, for a proof (the argument above for an $L^2$ convergence follows the one given in Nemirovski et al. [145]).

Under (H3), one can say much more on the asymptotic behavior of the algorithm by comparing it with an ordinary differential equation. The "ODE method," introduced in Ljung [120], is indeed a fundamental tool for the analysis of stochastic approximation algorithms. The correspondence between discrete and continuous

times is provided by the sequence $\alpha_t$. More precisely, let $\tau_0 = 0$ and $\tau_t = \tau_{t-1} + \alpha_t$, $t \geq 1$. From (H3), $\tau_t \to \infty$ when $t \to \infty$. Define the piecewise linear interpolation $x^\ell(\rho)$ of the sequence $x_t$ by

$$X^\ell(\rho) = X_t + \frac{\rho - \tau_t}{\alpha_{t+1}}(X_{t+1} - X_t), \quad \rho \in [\tau_t, \tau_{t+1}).$$

Switching to continuous time allows us to interpret the average iteration $\bar{x}_{t+1} = \bar{x}_t + \alpha_{t+1}\bar{H}(\bar{x}_t)$ as an Euler discretization scheme for the ordinary differential equation (ODE)

$$\partial_\rho \bar{x} = \bar{H}(\bar{x}). \tag{3.24}$$

Most of the insight on long-term behavior of stochastic approximations results from the fact that the random process $x$ behaves asymptotically like solutions of this ODE. One has, for example, the following result, for which we introduce some additional notation.

Assume that (3.24) has unique solutions for given initial conditions on any finite interval, and denote by $\varphi(\rho, \omega)$ its solution at time $\rho$ initialized with $\bar{x}(0) = \omega$. Let $\alpha^c(\rho)$ and $\eta^c(\rho)$ be piecewise constant interpolations of $(\alpha_t)$ and $(\eta_t)$ defined by $\alpha^c(\rho) = \alpha_{t+1}$ and $\eta^c(\rho) = \eta_{t+1}$ on the interval $[\tau_t, \tau_{t+1})$. Finally, let

$$\Delta(\rho, T) = \max_{s \in [\rho, \rho+T]} \left| \int_\rho^s \eta^c(u)du \right|.$$

The following proposition (see [23]) compares the tails of the process $x^\ell$ (i.e., the functions $x^\ell(\rho + s)$, $s \geq 0$) with the solutions of the ODE over finite intervals.

**Proposition 3.27 (Benaim)** *Assume that $\bar{H}$ is Lipschitz and bounded. Then, for some constant $C(T)$ that only depends on $T$ and $\bar{H}$, one has, for all $\rho \geq 0$*

$$\sup_{h \in [0,T]} |X^\ell(\rho + h) - \varphi(h, X^\ell(\rho))| \leq C(T)\left( \Delta(\rho - 1, T + 1) + \max_{s \in [\rho, \rho+T]} \alpha^c(s) \right). \tag{3.25}$$

Recall that $\bar{H}$ being Lipschitz means that there exists a constant $C$ such that

$$|\bar{H}(w) - \bar{H}(w')| \leq C|w - w'|$$

for all $w, w' \in \mathbb{R}^p$.

In the upper-bound in (3.25), the term $\Delta(\rho - 1, T + 1)$ is a random variable. It can be related to the variations

$$\Delta'(t, N) = \max_{k=0,\dots,N} |M_{t+k} - M_t|,$$

where $M$ is defined in (3.19), because, if $m(\rho)$ is the largest integer $t$ such that $\tau_t \leq \rho$, then

$$\Delta'(m(\rho)+1, m(\rho+T) - m(\rho)) \leq \Delta(\rho, T) \leq \Delta'(m(\rho), m(\rho+T) - m(t) + 1).$$

In the case we are considering, one can use martingale inequalities (called Doob's inequalities) to control $\Delta'$. One has, for example,

$$P\left(\max_{0 \leq k \leq N} |M_{t+k} - M_t| > \lambda\right) \leq \frac{E(|M_{t+N} - M_t|^2)}{\lambda^2}. \tag{3.26}$$

Furthermore, using the fact that $E(\eta_{k+1}\eta_{l+1}) = 0$ if $k \neq l$, one has

$$E(|M_{t+N} - M_t|^2) = \sum_{k=t}^{t+N} \alpha_{k+1}^2 E(|\eta_{t+1}|^2).$$

If we assume (to simplify) that $H$ is bounded and $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$ then, for some constant $C$, we have

$$E(|M_{t+N} - M_t|^2) \leq C \sum_{k=t}^{\infty} \alpha_{k+1}^2 \to 0$$

and inequality (3.26) can then be used in (3.25) to control the probability of deviation of the stochastic approximation from the solution of the ODE over finite intervals (a little more work is required under weaker assumptions on $H$, such as (H1)).

Proposition 3.27 cannot be used with $T = \infty$ because the constant $C(T)$ typically grows exponentially with $T$. In order to draw conclusions on the limit of the process $W$, one needs additional assumptions on the stability of the ODE. We refer to [23] for a collection of results on the relationship between invariant sets and attractors of the ODE and limit trajectories of the stochastic approximation. We here quote one of these results which is especially relevant for SGD.

**Proposition 3.28** *Assume that $\bar{H} = -\nabla E$ is the gradient of a function $E$ and that $\nabla E$ only vanishes at a finite number of points. Assume also that $X_t$ is bounded. Then $X_t$ converges to a point $x^*$ such that $\nabla E(x^*) = 0$.*

Some additional conditions on $\bar{H}$ can ensure that stochastic approximation trajectories remain bounded. The simplest one assumes the existence of a "Lyapunov function" that controls the ODE at infinity. The following result is a simplified version of Theorem 17 in Benveniste et al. [25].

**Theorem 3.29** *In addition to the hypotheses previously made, assume that there exists a $C^2$ function U with bounded second derivatives and $K_0 > 0$ such that, for all x such that $|x| \geq K_0$,*

$$\nabla U(x)^T \bar{H}(x) \leq 0,$$
$$U(x) \geq \gamma |x|^2, \gamma > 0.$$

*Then, the trajectories $X^\ell(\rho)$ are almost surely bounded.*

Note that hypothesis (H2) above implies the theorem's assumptions.

### 3.3.3   The ADAM algorithm

ADAM (for adaptive moment estimation [102]) is a popular variant of stochastic gradient descent. When dealing with high-dimensional vectors $\mathcal{W}$, using a single "gain" parameter ($\gamma_{n+1}$ in (11.2)) is a limiting assumption since all parameters do not need to scale the same way. This can sometimes be handled by reweighting the components of $H$, i.e., using iterations

$$X_{t+1} = X_t + \alpha_t D_t H(X_t, \xi_{t+1})$$

where $D_t$ is a (typically diagonal) matrix. The previous theory can be applied to situations in which $D$ may be random, provided it converges almost surely to a fixed matrix.

   The ADAM algorithm provides such a construction (without the theoretical guarantees) in which $D_t$ is computed using past iterations of the algorithm. It requires several parameters, namely: $\alpha$: the algorithm gain, taken as constant (e.g., $\alpha = 0.001$); Two parameters $\beta_1$ and $\beta_2$ for moment estimates (e.g. $\beta_1 = 0.9$ and $\beta_2 = 0.999$); A small number $\epsilon$ (e.g., $\epsilon = 10^{-8}$) to avoid divisions by 0. In addition, ADAM defines two vectors: a mean $m$ and a second moment $v$, respectively initialized at **0** and $\mathbb{1}$. The ADAM iterations are given below, in which $g^{\otimes 2}$ denotes the vector obtained by squaring each coefficient of a vector $g$.

---

**Algorithm 3.1 (ADAM)**
   1. Let $X_t$ be the current state, $m_t$ and $v_t$ the current mean and variance.

   2. Generate $\xi_{t+1}$ and let $g_{t+1} = H(X_t, \xi_{t+1})$.

   3. Update $m_{t+1} = \beta_1 m_t + (1 - \beta_1)g_{t+1}$.

   4. Update $v_{t+1} = \beta_2 v_t + (1 - \beta_2)g_{t+1}^{\odot 2}$.

   5. Let $\hat{m}_{t+1} = m_{t+1}/(1 - \beta_1^{t+1})$ and $\hat{v}_{t+1} = v_{t+1}/(1 - \beta_2^{t+1})$

6. Set

$$X_{t+1} = X_t - \alpha \frac{\hat{m}_{t+1}}{\sqrt{\hat{v}_{t+1}} + \epsilon}$$

---

Note that the iteration on $m_t$ and $v_t$ correspond to defining

$$\hat{m}_t = \frac{\beta_1}{1 - \beta_1^t} \sum_{k=0}^{t} (1 - \beta_1)^{t-k} g_k$$

and

$$\hat{v}_t = \frac{\beta_2}{1 - \beta_2^t} \sum_{k=0}^{t} (1 - \beta_2)^{t-k} g_k^{\odot 2}.$$

## 3.4 Constrained optimization problems

### 3.4.1 Lagrange multipliers

A constrained optimization problem minimizes a function $F$ over a closed subset $\Omega$ of $\mathbb{R}^d$, with $\Omega \neq \mathbb{R}^d$. This restriction invalidates, in a large part, the optimality conditions discussed in section 3.2. These conditions indeed apply to minimizers belonging to the interior of $\Omega$, and therefore do not hold when they lie at its boundary, which is a very common situation in practice ($\Omega$ often has an empty interior).

In this section, which follows the discussion given in Wright and Recht [204], we review conditions for optimality for constrained minimization of smooth functions, in two cases. The first one, discussed in this section, is when $\Omega$ is defined by a finite number of smooth constraints, leading, under some assumptions, to the Karush-Kuhn-Tucker (or KKT) conditions. The second one, in the next section, specializes to closed convex $\Omega$.

**KKT conditions**

We introduce some notation. Let $\gamma_i$, for $i \in \mathcal{C}$, be $C^1$ functions $\gamma_i : \mathbb{R}^d \to \mathbb{R}$, where $\mathcal{C}$ is a finite set of indices. We assume that $\mathcal{C}$ is divided into two non-intersecting parts, $\mathcal{C} = \mathcal{E} \cup \mathcal{I}$ and consider minimization problems searching for

$$x^* \in \operatorname*{argmin}_{\Omega} F \tag{3.27}$$

where

$$\Omega = \{x \in \mathbb{R}^d : \gamma_i(x) = 0, i \in \mathcal{E} \text{ and } \gamma_i(x) \leq 0, i \in \mathcal{I}\}. \tag{3.28}$$

The set $\Omega$ of all $x$ that satisfy the constraints is called the *feasible set* for the considered problem. We will always assume that it is non-empty. If $x \in \Omega$, one defines the set $\mathcal{A}(x)$ of *active constraints* at $x$ to be

$$\mathcal{A}(x) = \{i \in \mathcal{C} : \gamma_i(x) = 0\}.$$

One obviously has $\mathcal{E} \subset \mathcal{A}(x)$ for $x \in \Omega$.

To be valid, the KKT conditions require some additional assumptions on potential minimizers, called "constraint qualifications." An instance of such assumptions is provided by the next definition.

**Definition 3.30** *A point $x \in \Omega$ satisfies the Mangasarian-Fromovitz constraint qualifications (MF-CQ) if the following two conditions are satisfied.*

*(MF1)  The vectors $(\nabla \gamma_i(x), i \in \mathcal{E})$ are linearly independent.*

*(MF2)  There exists a vector $h \in \mathbb{R}^d$ such that $h^T \nabla \gamma_i(x) = 0$ for all $i \in \mathcal{E}$ and $h^T \nabla \gamma_i(x) < 0$ for all $i \in \mathcal{A}(x) \cap \mathcal{I}$.*

A sufficient (and easier to check) condition for $x$ to satisfy these constraints is when the vectors $(\nabla \gamma_i(x), i \in \mathcal{A}(x))$ are linearly independent [37]. Indeed, if the latter "LI-CQ" condition is true, then any set of values can be assigned to $h^T \nabla \gamma_i(x)$ with the existence of a vector $h$ that achieves them.

We introduce the *Lagrangian*

$$L(x, \lambda) = F(x) + \sum_{i \in \mathcal{C}} \lambda_i \gamma_i(x) \tag{3.29}$$

where the real numbers $\lambda_i$, $i \in \mathcal{C}$ are called *Lagrange multipliers*. The following theorem (stated without proof, see, e.g., [146, 35]) provides necessary conditions satisfied by solutions of the constrained minimization problem that satisfy the constraint qualifications.

**Theorem 3.31** *Assume $x^* \in \Omega$ is a solution of (3.27), and that $x^*$ satisfies the MF-CQ conditions. Then there exist Lagrange multipliers $\lambda_i$, $i \in \mathcal{C}$, such that*

$$\begin{cases} \partial_x L(x^*, \lambda) = 0 \\ \lambda_i \geq 0 \text{ if } i \in \mathcal{I}, \text{ with } \lambda_i = 0 \text{ when } i \notin \mathcal{A}(x^*) \end{cases} \tag{3.30}$$

Conditions (3.30) are the KKT conditions for the constrained optimization problem. The second set of conditions is often called the *complementary slackness conditions* and state that $\lambda_i = 0$ for an inequality constraint unless this constraint is satisfied with an equality. The next section provides examples in which the MF-CQ conditions are not satisfied and Theorem 3.31 does not hold. However, these conditions are not needed in the special case when the constraints are affine.

**Theorem 3.32** *Assume that for all $i \in \mathcal{A}(x^*)$, the functions $\gamma_i$ are affine, i.e., $\gamma_i(x) = b_i^T x + \beta_i$ for some $b \in \mathbb{R}^d$ and $\beta \in \mathbb{R}$. Then (3.30) holds at any solution of (3.27).*

**Remark 3.33** We have taken the convention to express the inequality constraints as $\gamma_i(x) \leq 0$, $i \in \mathcal{I}$. With the reverse convention, i.e., $\gamma_i(x) \geq 0$, $i \in \mathcal{I}$, one generally defines the Lagrangian as

$$L(x, \lambda) = F(x) - \sum_{i \in \mathcal{C}} \lambda_i \gamma_i(x)$$

and the KKT conditions remain unchanged. ♦

**Examples.** Constraint qualifications are important to ensure the validity of the theorem. Consider a problem with equality constraints only, and replace it by

$$x^* \in \underset{\Omega}{\operatorname{argmin}} F$$

subject to $\tilde{\gamma}_i(x) = 0$, $i \in \mathcal{E}$, with $\tilde{\gamma}_i = \gamma_i^2$. We clearly did not change the problem. However, the previous theorem applied to the Lagrangian

$$L(x, \lambda) = F(x) + \sum_{i \in \mathcal{C}} \lambda_i \tilde{\gamma}_i(x)$$

would require an optimal solution to satisfy $\nabla F(x) = 0$, because $\nabla \tilde{\gamma}_i(x) = 2\gamma_i(x)\nabla\gamma_i(x) = 0$ for any feasible solution. Minimizers of constrained problems do not necessarily satisfy $\nabla F(x) = 0$, however. This is no contradiction with the theorem since $\nabla \tilde{\gamma}_i(x) = 0$ for all $i$ shows that no feasible point satisfies the MF-CQ.

To take a more specific example, still with equality constraints, let $d = 3, \mathcal{C} = \{1, 2\}$ with $F(x, y, z) = x/2 + y$ and $\gamma_1(x, y, z) = x^2 - y^2, \gamma_2(x, y, z) = y - z^2$. Note that $\gamma_1 = \gamma_2 = 0$ implies that $y = |x|$, so that, for a feasible point, $F(x, y, z) = |x| + x/2 \geq 0$ and vanishes only when $x = y = 0$, in which case $z = 0$ also. So $(0, 0, 0)$ is a global minimizer. We have $dF(0) = (1/2, 1, 0)$, $d\gamma_1(0) = (0, 0, 0)$ and $d\gamma_2(0) = (0, 1, 0)$ so that $0$ does not satisfy the MF-CQ. The equation

$$dF(0) + \lambda_1 d\gamma_1(0) + \lambda_2 d\gamma_2(0) = 0$$

has no solution $(\lambda_1, \lambda_2)$, so that the conclusion of the theorem does not hold.

### 3.4.2 Convex constraints

We now consider the case in which $\Omega$ is a closed convex set. To specify the optimality conditions in this case, we need the following definition.

**Definition 3.34** *Let $\Omega \subset \mathbb{R}^d$ be convex and let $x \in \Omega$. The normal cone to $\Omega$ at $x$ is the set*

$$\mathcal{N}_\Omega(x) = \{h \in \mathbb{R}^d : h^T(y - x) \le 0 \text{ for all } y \in \Omega\} \tag{3.31}$$

The normal cone is an example of convex cone. (A convex subset $\Gamma$ of $\mathbb{R}^d$ is called a convex cone, if it is such that $\lambda x \in \Gamma$ for all $x \in \Gamma$ and $\lambda \ge 0$, a property obviously satisfied by $\mathcal{N}_\Omega(x)$.) It should also be clear from the definition that non-zero vectors in $\mathcal{N}_\Omega(x)$ always point outside $\Omega$, i.e., $x + h \notin \Omega$ if $h \in \mathcal{N}_\Omega(x)$, $h \ne 0$. Here are some examples.

- If $x$ is in the interior of $\Omega$, then $\mathcal{N}_\Omega(x) = \{0\}$.

- Assume that $\Omega$ is a half space, i.e., $\Omega = \{x : b^T x + \beta \le 0\}$ with $|b| = 1$, and take $x \in \partial\Omega$, i.e., $b^T x + \beta = 0$. Then

$$\mathcal{N}_\Omega(x) = \{h = \mu b : \mu \ge 0\}.$$

Indeed, any element of $\mathbb{R}^d$ can be written as $y = x + \lambda b + q$ with $q^T b = 0$, and $y \in \Omega$ if and only if $\lambda \le 0$. Fix such a $y$ and take $h \in \mathbb{R}^d$, decomposed as $h = \mu b + r$, with $r^T b = 0$. We have $h^T(y - x) = \lambda\mu + r^T q$. Clearly, if $\mu < 0$, or if $r \ne 0$, one can find $\lambda \le 0$ and $q \perp b$ such that $h^T(y - x) > 0$. One the other hand, if $\mu \le 0$ and $r = 0$, we have $h^T(y - x) \le 0$ for all $y \in \Omega$, which proves the above statement.

- With a similar argument, if $\Omega = \{x : b^T x + \beta = 0\}$ is a hyperplane, one finds that

$$\mathcal{N}_\Omega(x) = \{h = \lambda b : \lambda \in \mathbb{R}\}.$$

One can build normal cones to domains associated with multiple inequalities or equalities based on the following theorem.

**Theorem 3.35** *Let $\Omega_1$ and $\Omega_2$ be two convex sets with $\text{relint}(\Omega_1) \cap \text{relint}(\Omega_2) \ne \emptyset$. Then, if $x \in \Omega_1 \cap \Omega_2$*

$$\mathcal{N}_{\Omega_1 \cap \Omega_2}(x) = \mathcal{N}_{\Omega_1}(x) + \mathcal{N}_{\Omega_2}(x)$$

*Here, the addition is the standard sum between sets in a vector space:*

$$A + B = \{x + y : x \in A, y \in B\}.$$

Finally, we note that, if $x \in \text{relint}(\Omega)$, then

$$\mathcal{N}_\Omega(x) = \{h \in \mathbb{R}^d : h^T(y - x) = 0, y \in \Omega\}. \tag{3.32}$$

Indeed, if $y \in \Omega$, then $x + \epsilon(y - x) \in \Omega$ for small enough $\epsilon$ (positive or negative). For $h \in \mathcal{N}_\Omega(x)$, the condition $\epsilon h^T(y - x) \le 0$ for small enough $\epsilon$ requires that $h^T(y - x) = 0$.

With this definition in hand, we have the following theorem.

**Theorem 3.36** *Let F be a $C^1$ function and $\Omega$ a closed convex set. If $x^* \in \mathrm{argmin}_\Omega F$, then*

$$-\nabla F(x^*) \in \mathcal{N}_\Omega(x^*). \tag{3.33}$$

*If F is convex and (3.33) holds, then $x^* \in \mathrm{argmin}_\Omega F$.*

Proof Assume that $x^* \in \mathrm{argmin}_\Omega F$. If $y \in \Omega$, then $x^* + t(y - x^*) \in \Omega$ for all $t \in [0, 1]$ and the function $f(t) = F(x + t(y - x^*))$ is $C^1$ on $[0, 1]$, with a minimum at $t = 0$. This requires that $\partial_t f(0) = \nabla F(x^*)^T (y - x^*) \geq 0$, because, if $\partial_t f(0) < 0$, a Taylor expansion would show that $f(t) < f(0)$ for small enough $t > 0$.

If $F$ is convex and (3.33) holds, we have $F(y) \geq F(x^*) + \nabla F(x^*)^T (y - x^*)$ by convexity, so that
$$F(x^*) \leq F(y) + (-\nabla F(x^*))^T (y - x^*) \leq F(y). \qquad \blacksquare$$

### 3.4.3 Applications

**Lagrange multipliers revisited.** Consider $\Omega$ defined by (3.28), with the additional assumptions that $\gamma_i(x) = b_i^T x + \beta_i$ for $i \in \mathcal{E}$ and $\gamma_i$ is convex for $i \in \mathcal{I}$, which ensure that $\Omega$ is convex. Define

$$\mathcal{N}'_\gamma(x) = \left\{ g = \sum_{i \in \mathcal{A}(x)} \lambda_i \nabla \gamma_i(x) : \lambda_i \geq 0, i \in \mathcal{A}(x) \cap \mathcal{I} \right\}.$$

Then, the KKT conditions in (3.30) can be rewritten as

$$-\nabla F(x^*) \in \mathcal{N}'_\gamma(x^*).$$

Note that one always have $\mathcal{N}'_\gamma(x) \subset \mathcal{N}_\Omega(x)$ since, for $g = \sum_{i \in \mathcal{A}(x)} \lambda_i \nabla \gamma_i(x) \in \mathcal{N}'_\gamma(x)$, one has, for $y \in \Omega$,

$$
\begin{aligned}
g^T(y - x) &= \sum_{i \in \mathcal{A}(x)} \lambda_i \nabla \gamma_i(x)^T (y - x) \\
&= \sum_{i \in \mathcal{E}} \lambda_i (a_i^T y - a_i^T x) + \sum_{i \in \mathcal{A}(x) \cap \mathcal{I}} \lambda_i (\gamma_i(x) + \nabla \gamma_i(x)^T (y - x)) \\
&= \sum_{i \in \mathcal{A}(x) \cap \mathcal{I}} \lambda_i (\gamma_i(x) + \nabla \gamma_i(x)^T (y - x)) \\
&\leq \lambda_i \gamma_i(y) \leq 0,
\end{aligned}
$$

in which the have used the facts that $a_i^T x = a_i^T y = -\beta_i$ for $x, y \in \Omega$, $i \in \mathcal{E}$, $\gamma_i(x) = 0$ for $i \in \mathcal{A}(x)$ and the convexity of $\gamma_i$. Constraint qualifications such as those considered above are sufficient conditions that ensure the identity between the two sets.

Consider now the situation of theorem 3.32, and assume that all constraints are affine inequalities, $\gamma_i(x) = b_i^T x + \beta \le 0, i \in \mathcal{I}$. Then, the statement $\mathcal{N}_\Omega(x) \subset \mathcal{N}_\gamma'(x)$ can be reexpressed as follows. All $h \in \mathbb{R}^d$ such that

$$h^T(y - x) \le 0$$

as soon as $b_i^T(y - x) \le 0$ for all $i \in \mathcal{A}(x)$ must take the form

$$h = \sum_{i \in \mathcal{A}(x)} \lambda_i b_i$$

with $\lambda^{(i)} \ge 0$. This property is called *Farkas's lemma* (see, e.g. [167]). Note that affine equalities $b_i^T x + \beta = 0$ can be included as two inequalities $b_i^T x + \beta \le 0$, $-b_i^T x - \beta \le 0$, which removes the sign constraint on the corresponding $\lambda^{(i)}$ and therefore yields theorem 3.32.

**Positive semi-definite matrices.** We now take an example in which theorem 3.32 does not apply directly. Let $\Omega = \mathcal{S}_n^+$ be the space of positive semidefinite $n \times n$ matrices, considered as a subset of the space $\mathcal{M}_n$ of $n \times n$ matrices, itself identified with $\mathbb{R}^{n^2}$. With this identification, the Euclidean inner product between two matrices can be expressed as $(A, B) \mapsto \mathrm{trace}(A^T B)$.

We have $A \in \mathcal{S}_n^+$ if and only if, for all $u \in \mathbb{R}^d$, $u^T A u \ge 0$, which provides an infinity of linear inequality constraints on $A$. Elements of $\mathcal{N}_{\mathcal{S}^+}(A)$ are matrices $H \in \mathcal{M}_n$ such that

$$\mathrm{trace}(H^T(B - A)) \le 0$$

for all $B \in \mathcal{S}_n^+$, and we want to make this normal cone explicit. We first note that, every square matrix $H$ can be decomposed as the sum of a symmetric matrix, $H_s$ and of a skew symmetric one, $H_a$ (namely, $H_s = (H + H^T)/2$ and $H_a = (H - H^T)/2$). We have moreover $\mathrm{trace}(H_a^T(B - A)) = 0$, so the condition is only on the symmetric part of $H$.

For any $u \in \mathbb{R}^d$, one can take $B = A + uu^T$, which belongs to $\mathcal{S}_n^+$, with $\mathrm{trace}(H_s^T(B - A)) = u^T H_s u$. This shows that, for $H$ to belong to $\mathcal{N}_{\mathcal{S}_n^+}(A)$, one needs $H_s \le 0$.

Now, take an eigenvector $u$ of $A$ with eigenvalue $\rho > 0$. Then $B = A - \alpha uu^T$ is also in $\mathcal{S}_n^+$ as soon as $0 \le \alpha \le \rho$, and $\mathrm{trace}(H_s^T(B - A)) = -\alpha u^T H_s u$. So, if $H \in \mathcal{N}_{\mathcal{S}_n^+}(A)$, we have $u^T H_s u \ge 0$, and since $H_s \le 0$, this gives $u^T H_s u = 0$. Still because $H_s$ is negative semi-definite, this implies $H_s u = 0$. (This can be shown, for example, using Schwarz's inequality which says that $(u^T H_s v)^2 \le (u^T H_s u)(v^T H_s v)$ for all $v \in \mathbb{R}^d$.) Decomposing $A$ with respect to its non-zero eigenvectors, i.e., writing

$$A = \sum_{k=1}^p \rho_k u_k u_k^T$$

where $p = \text{rank}(A)$, we get $AH_s = H_sA = 0$. We therefore obtained the proposition

**Proposition 3.37** *Let $A \in \mathcal{S}_n^+$. Then $H \in \mathcal{M}_n$ belongs to $\mathcal{N}_{\mathcal{S}_n^+}(A)$ if and only if $-H_s \in \mathcal{S}_n^+$ and $H_sA = 0$, where $H_s = (H + H^T)/2$.*

Now, if one wants to minimize a function $F$ over positive semidefinite matrices, and $A^*$ is a minimizer, we get the necessary condition that $A^*(\nabla F(A^*))_s = 0$ with $(\nabla F(A^*))_s$ positive semidefinite. These conditions are sufficient if $F$ is convex.

For example, take

$$F(A) = \frac{1}{2}\text{trace}(A^2) - \text{trace}(BA) \tag{3.34}$$

with $B \in \mathcal{S}_n$. Then $(\nabla F(A))_s = A - B$ and the condition is $A(A - B) = 0$ with $A \succeq B$. If $B$ is diagonalized in the form $B = U^T D U$, with $U$ orthogonal and $D$ diagonal, then the solution is $A^* = U^T D^+ U$ where $D^+$ is deduced from $U$ by replacing non-negative entries by zeros.

**Projection.** Let $\Omega$ be closed convex, $x_0 \in \mathbb{R}^d$ and $F(x) = \frac{1}{2}|x - x_0|^2$. We have

$$\min_\Omega F = \min_{\Omega \cap \bar{B}(0,R)} F$$

for large enough $R$ (e.g., larger than $F(x)$ for any fixed point in $\Omega$), and since the latter minimization is over a compact set, $\text{argmin}_\Omega F$ is not empty. The function $F$ being strongly convex, its minimizer over $\Omega$ is unique and called the projection of $x_0$ on $\Omega$, denoted $\text{proj}_\Omega(x_0)$.

Since $\nabla F(x) = x - x_0$, theorem 3.36 implies that $\text{proj}_\Omega(x_0)$ is characterized by $\text{proj}_\Omega(x_0) \in \Omega$ and

$$x_0 - \text{proj}_\Omega(x_0) \in \mathcal{N}_\Omega(\text{proj}_\Omega(x_0)) \tag{3.35}$$

or

$$(x_0 - \text{proj}_\Omega(x_0))^T(y - \text{proj}_\Omega(x_0)) \leq 0 \text{ for all } y \in \Omega. \tag{3.36}$$

If $x_0 \notin \Omega$, then $\text{proj}_\Omega(x_0) \in \partial\Omega$, since otherwise we would have $\mathcal{N}_\Omega(\text{proj}_\Omega(x_0)) = \{0\}$ and $x_0 = \text{proj}_\Omega(x_0)$, a contradiction. Of course, if $x_0 \in \Omega$, then $\text{proj}_\Omega(x_0) = x_0$.

Here are some important examples.

**1.** Let $\Omega = z_0 + V$, where $z_0 \in \mathbb{R}^d$ and $V$ is a linear space (i.e., $\Omega$ is an affine subset of $\mathbb{R}^d$). Then $N_\Omega(x) = z_0 + V^\perp = x + V^\perp$ for all $x \in \Omega$, where $V^\perp$ is the vector space of vectors orthogonal to $V$, and $\text{proj}_\Omega(x_0)$ is characterized by $\text{proj}_\Omega(x_0) \in \Omega$ and

$$(x_0 - \text{proj}_\Omega(x_0)) \in V^\perp$$

which is the usual characterization of the orthogonal projection on an affine space (compare to section 6.4).

**2.** If $\Omega = \bar{B}(0, 1)$, the closed unit sphere, then $N_\Omega(x) = \mathbb{R}^+ x$ for $x \in \partial\Omega$ (i.e., $|x| = 1$). One can indeed note that, if $h \neq 0$ in normal to $\Omega$ at $x$, then $h/|h| \in \Omega$ so that

$$h^T \left( \frac{h}{|h|} - x \right) \leq 0$$

which yields $|h| \leq h^T x$. The Cauchy-Schwartz inequality implying that $h^T x \leq |h||x| = |h|$, we must have equality, $h^T x = |h||x|$, which is only possible when $x$ and $h$ are collinear.

Given $x_0 \in \mathbb{R}^d$ with $x_0 \geq 1$, we see that $\text{proj}_\Omega(x_0)$ must satisfy the conditions $|\text{proj}_\Omega(x_0)| = 1$ (to be in $\partial\Omega$) and $x_0 - \text{proj}_\Omega(x_0) = \lambda x_0$ for some $\lambda \geq 0$, which gives $\text{proj}_\Omega(x_0) = x_0/|x_0|$.

**3.** If $\Omega = \mathcal{S}_n^+$ and $B$ (taking the role of $x_0$) is a symmetric matrix, then $\text{proj}_\Omega(B)$ was found in the previous section, and is given by $A = U^T D^+ U$ where $U^T D U$ provides a diagonalization of $B$.

The projection has the important property of being 1-Lipschitz.

**Proposition 3.38** *Let $\Omega$ be a closed convex subset of $\mathbb{R}^d$. Then, for all $x, y \in \mathbb{R}^d$*

$$|\text{proj}_\Omega(x) - \text{proj}_\Omega(y)| \leq |x - y|. \tag{3.37}$$

PROOF This proposition is a special case of proposition 3.55 below.                    ∎

### 3.4.4   Projected gradient descent

The projected gradient descent algorithm minimizes $F$ over $\Omega$ by iterating

$$x_{t+1} = \text{proj}_\Omega(x_t - \alpha_t \nabla F(x_t)), \tag{3.38}$$

which provides a feasible method when $\text{proj}_\Omega$ is easy to compute. An equivalent formulation is

$$x_{t+1} = \underset{\Omega}{\text{argmin}} \, F(x_t) + \nabla F(x_t)^T (x - x_t) + \frac{1}{2\alpha_t}|x - x_t|^2. \tag{3.39}$$

To justify this last statement it suffices to notice that the function in the r.h.s. can be written as

$$\frac{1}{2\alpha_t}|x - x_t + \alpha_t \nabla F(x_t)|^2 - \frac{\alpha_t}{2}|\nabla F(x_t)|^2 + F(x_t)$$

and apply the definition of the projection.

The convergence properties of this algorithm will be discussed in section 3.5.5, in a more general context.

## 3.5 General convex problems

### 3.5.1 Epigraphs

**Definition 3.39** *Let F be a convex function. The epigraph of F is the set*

$$\text{epi}(F) = \left\{(x, a) \in \mathbb{R}^d \times \mathbb{R} : F(x) \le a\right\}. \tag{3.40}$$

*One says that F is closed if* $\text{epi}(F)$ *is a closed subset of* $\mathbb{R}^d \times \mathbb{R}$, *that is: if* $x = \lim_n x_n$ *and* $a = \lim_n a_n$ *with* $F(x_n) \le a_n$, *then* $F(x) \le a$.

Clearly, if $(x, a) \in \text{epi}(F)$, then $x \in \text{dom}(F)$. It should also be clear that $\text{epi}(F)$ is always convex when $F$ is convex: If $(x, a), (y, b) \in \text{epi}(F)$, then

$$F((1 - t)x + ty) \le (1 - t)F(x) + tF(y) \le (1 - t)a + tb$$

so that $(1 - t)(x, a) + t(y, b) \in \text{epi}(F)$.

To illustrate the definition, consider a simple example. Let $F$ be the function defined on $\mathbb{R}$ by $F(x) = |x|$ if $|x| < 1$ and $F(x) = +\infty$ otherwise. It is convex, but not closed, as can be seen by taking the sequence $(1 - 1/n, 1) \in \text{epi}(F)$, with, at the limit, $F(1) = +\infty > 1$. In contrast, the function defined by $\tilde{F}(x) = |x|$ if $|x| \le 1$ and $\tilde{F}(x) = +\infty$ otherwise is convex and closed.

We have the following proposition.

**Proposition 3.40** *A convex function F is closed if and only if all its sub-level sets*

$$\Lambda_a(F) = \left\{x \in \mathbb{R}^d : F(x) \le a\right\}$$

*are closed subsets of* $\mathbb{R}^d$.

PROOF If $F$ is closed, then $\Lambda_a(F)$ is the intersection of the set $\{(x, a) : x \in \mathbb{R}^d\}$, which is obviously closed, and of $\text{epi}(F)$. It is therefore a closed set.

Conversely, assume that all $\Lambda_a(F)$ are closed and take a sequence $(x_n, a_n)$ in $\text{epi}(F)$ that converges to $(x, a)$. Then, fixing $\epsilon > 0$, $x_n \in \Lambda_{a+\epsilon}$ for large enough $n$, and since this set is closed, $F(x) \le a + \epsilon$. Since this is true for all $\epsilon > 0$, we have $F(x) \le a$ and $(x, a) \in \text{epi}(F)$. $\blacksquare$

Note that, if $F$ is continuous, then it is closed, so that closedness generalizes continuity for convex functions, but it also applies to the non-smooth case.

If $\Omega$ is a convex subset of $\mathbb{R}^d$, its indicator function $\sigma_\Omega$ (such that $\sigma_\Omega(x) = 0$ for $x \in \Omega$ and $\sigma_\Omega(x) = +\infty$ otherwise) is closed if and only if $\Omega$ is a closed subset of $\mathbb{R}^d$. This is obvious since $\Lambda_a(\sigma_\Omega) = \Omega$ if $a \ge 0$ and $\emptyset$ otherwise.

### 3.5.2   Subgradients

Several machine learning problems involve convex functions that are not $C^1$, requiring a generalization of the notion of derivative provided by the following definition.

**Definition 3.41** *If F is a convex function and $x \in \mathrm{dom}(F)$, a vector $g \in \mathbb{R}^d$ such that*

$$F(x) + g^T(y - x) \leq F(y) \tag{3.41}$$

*for all $y \in \mathbb{R}^d$ is called a* subgradient *of F at x.*

*The set of subgradients of F at x is denoted $\partial F(x)$ and called the* subdifferential *of F at x.*

If $x \in \mathrm{int}(\mathrm{dom}(F))$ and $F$ is differentiable at $x$, (3.5) implies that $\nabla F(x) \in \partial F(x)$. This is in this case the only element of $\partial F(x)$.

**Proposition 3.42** *If F is differentiable at $x \in \mathrm{int}(\mathrm{dom}(F))$, then $\partial F(x) = \{\nabla F(x)\}$.*

PROOF  We need to prove that there is no other subgradient. Assume that $\nabla F(x)$ exists and take $y = x + \epsilon u$ in (3.41) ($u \in \mathbb{R}^d$). One gets, for $g \in \partial F(x)$,

$$\epsilon g^T u \leq F(x + \epsilon u) - F(x) = \epsilon \nabla F(x)^T u + o(\epsilon)$$

Dividing by $|\epsilon|$ and letting $\epsilon \to 0$ gives (depending on the sign of $\epsilon$)

$$g^T u \leq \nabla F(x)^T u \text{ and } -g^T u \leq -\nabla F(x)^T u$$

This is only possible if $g^T u = \nabla F(x)^T u$ for all $u \in \mathbb{R}^d$, which itself implies $g = \nabla F$.  ∎

The next theorem, which is an obvious consequence of definition 3.41, characterizes minimizers of convex functions in the general case.

**Theorem 3.43** *Let $F : \mathbb{R}^d \to \mathbb{R}$ be convex. Then x is a (global) minimizer of F if and only if $0 \in \partial F(x)$.*

The following result shows that subgradients exist under generic conditions. We note that $g \in \partial F(x)$ if and only if $\mathrm{proj}_{\overrightarrow{\mathrm{aff}(\mathrm{dom}(F))}}(g) \in \partial F$, because (3.41) is trivial if $F(y) = +\infty$. So $\partial F$ cannot be bounded unless $\mathrm{aff}(\mathrm{dom}(D)) = \mathbb{R}^d$. However, it is the part of this set that is included in the $\overrightarrow{\mathrm{aff}}(\mathrm{dom}(F))$ that is of interest.

**Proposition 3.44** *For all $x \in \mathbb{R}^d$, $\partial F(x)$ is a closed convex set (possibly empty, in particular for $x \notin \mathrm{dom}(F)$). If $x \in \mathrm{ridom}(F)$, then $\partial F(x) \neq \emptyset$ and $\partial F(x) \cap \overrightarrow{\mathrm{aff}}(\mathrm{dom}(F))$ is compact.*

PROOF The convexity and closedness of $\partial F(x)$ is clear from the definition. If $x \in$ ridom($F$), there exists $\epsilon > 0$ such that $x + \epsilon h \in$ ridom($F$) for all $h \in \overrightarrow{\text{aff}}(\text{dom}(F))$ with $|h| = 1$. For all $g \in \partial F(x) \cap \overrightarrow{\text{aff}}(\text{dom}(F))$, one has

$$|g| = \max\{g^T h : h \in \overrightarrow{\text{aff}}(\text{dom}(F)), |h| = 1\}$$
$$\leq \max((F(x + \epsilon h) - F(x))/\epsilon : h \in \overrightarrow{\text{aff}}(\text{dom}(F)), |h| = 1)$$

and the upper bound is finite because it is the maximum of a continuous function over a bounded set. This shows that $\partial F(x)$ is bounded. We defer the proof that $\partial F(x) \neq \emptyset$ to section 3.7. ∎

Subdifferentials are, under mild conditions, additive. More precisely, we have the following proposition.

**Theorem 3.45** *Let $F_1$ and $F_2$ be convex functions such that*

$$\text{ridom}(F_1) \cap \text{ridom}(F_2) \neq \emptyset.$$

*Then, for all $x \in \mathbb{R}^d$, $\partial(F_1 + F_2)(x) = \partial F_1(x) + \partial F_2(x)$.*

Note that the inclusion

$$\partial F_1(x) + \partial F_2(x) \subset \partial(F_1 + F_2)(x)$$

as can be immediately checked by summing the inequalities satisfied by subgradients. The reverse inclusion requires the use of separation theorems for convex sets (see section 3.7).

Another important point is how the chain rule works with compositions with affine functions.

**Theorem 3.46** *Let $F$ be a convex function on $\mathbb{R}^d$, $A$ a $d \times m$ matrix and $b \in \mathbb{R}^d$. Let $G(x) = F(Ax + b)$, $x \in \mathbb{R}^m$. Assume that there exists $x_0 \in \mathbb{R}^m$ such that $Ax_0 \in$ ridom($F$). Then, for all $x \in \mathbb{R}^m$,*

$$\partial G(x) = A^T \partial F(Ax + b).$$

One direction is straightforward and does not require the condition on ridom($F$). If $g \in \partial F(Ax + b)$, then

$$F(z) - F(Ax + b) \geq g^T(z - Ax - b), z \in \mathbb{R}^d$$

and applying this inequality to $z = Ay + b$ for $y \in \mathbb{R}^m$ yields

$$G(y) - G(x) \geq g^T A(y - x)$$

so that $A^T g \in \partial G$ and $A^T \partial F \subset \partial G$. The reverse inclusion is proved in section 3.7.

Subdifferentials can be seen as generalizations of normal cones.

**Proposition 3.47** *Assume that $\Omega$ is a closed convex subset of $\mathbb{R}^d$. Then $\sigma_\Omega$ (the indicator function of $\Omega$) has a subdifferential everywhere on $\Omega$ with*

$$\partial\sigma_\Omega(x) = \mathcal{N}_\Omega(x), \ x \in \Omega$$

PROOF For $x \in \Omega$, (3.41) is

$$g^T(y - x) \le \sigma_\Omega(y)$$

for $y \in \mathbb{R}^d$, but since $\sigma_\Omega(y) = +\infty$ outside of $\Omega$, $g \in \partial\sigma_\Omega(x)$ is equivalent to

$$g^T(y - x) \le 0$$

for $y \in \Omega$, which is exactly the definition of the normal cone. ∎

Given this proposition, it is also clear (after noting that $\sigma_{\Omega_1} + \sigma_{\Omega_2} = \sigma_{\Omega_1 \cap \Omega_2}$) that theorem 3.45 is a generalization of theorem 3.35.

### 3.5.3 Directional derivatives

From proposition 3.5, applied with $y = x + h$, we see that

$$t \mapsto \frac{1}{t}(F(x + th) - F(x))$$

is increasing as a function of $t$. This property allows us to define directional derivatives of $F$ at $x$.

**Definition 3.48** *Let $F$ be convex and $x \in \mathrm{dom}(F)$. The directional derivative of $F$ at $x$ in the direction $h \in \mathbb{R}^d$ is defined by*

$$dF(x, h) = \lim_{t \downarrow 0} \frac{1}{t}(F(x + th) - F(x)), \tag{3.42}$$

*and belong to $[-\infty, +\infty]$.*

Note that, still from proposition 3.5, one has, for all $x \in \mathrm{dom}(F)$ and $y \in \mathbb{R}^d$:

$$F(y) \ge F(x) + dF(x, y - x) \tag{3.43}$$

We have the proposition:

**Proposition 3.49** *If $F$ is convex, then $x^* \in \mathrm{argmin}(F)$ if and only if $dF(x^*, h) \ge 0$ for all $h \in \mathbb{R}^d$.*

PROOF If $dF(x^*, h) \geq 0$, then $F(x^* + th) - F(x^*) \geq 0$ for all $t > 0$ and this being true for all $h$ implies that $x^*$ is a minimizer. Conversely, if $x^*$ is a minimizer, $dF(x^*, h)$ is a limit of non-negative numbers and is therefore non-negative. ∎

**Proposition 3.50** *If F is convex and $x \in \mathrm{dom}(F)$, then $dF(x, h)$ is positively homogeneous and subadditive (hence convex) as a function of h, namely*

$$dF(x, \lambda h) = \lambda dF(x, h), \lambda > 0$$

*and*

$$dF(x, h_1 + h_2) \leq dF(x, h_1) + dF(x, h_2).$$

PROOF Positive homogeneity is straightforward and left to the reader. For the second one, we can write

$$F(x + th_1 + th_2) \leq \frac{1}{2}(F(x + th_1/2) + F(x + th_2/2))$$

by convexity so that

$$\frac{1}{t}(F(x + th_1 + th_2) - F(x)) \leq \frac{1}{2}\left(\frac{1}{t}(F(x + th_1/2) - F(x)) + \frac{1}{t}(F(x + th_2/2) - F(x))\right).$$

Taking $t \downarrow 0$,

$$dF(x; h_1 + h_2) \leq \frac{1}{2}(dF(x; h_1/2) + dF(x, h_2/2)) = dF(x, h_1) + dF(x, h_2).$$ ∎

**Proposition 3.51** *If F is convex and $x \in \mathrm{dom}(F)$, then*

$$dF(x, h) \geq \sup\{g^T h, g \in \partial F(x)\}.$$

*If $x \in \mathrm{ridom}(F)$, then*

$$dF(x, h) = \max\{g^T h, g \in \partial F(x)\}.$$

PROOF If $g \in \partial F(x)$, then for all $t > 0$

$$F(x + th) - F(x) \geq tg^T h.$$

Dividing by $t$ and passing to the limit yields $dF(x, h) \geq g^T h$.

We prove that the maximum is attained at some $g \in \partial F(x)$ when $x \in \mathrm{ridom}(F)$. In this case, the domain of the convex function $G : \tilde{h} \mapsto dF(x, \tilde{h})$ is the vector space parallel to $\mathrm{aff}(\mathrm{dom}(F))$, namely

$$\mathrm{dom}(G) = \{h : x + h \in \mathrm{aff}(\mathrm{dom}(F))\}.$$

Indeed, for any $h$ in this set, there exists $\epsilon > 0$ such that $x + th \in \text{dom}(F)$ for $0 < t < \epsilon$ and $dF(x, h) \leq (F(x + th) - F(x))/t < \infty$. Conversely, if $h \in \text{dom}(G)$, then $F(x + th) - F(x)$ must be finite for small enough $t$, so that $x + th \in \text{dom}(F)$ and $x + h \in \text{aff}(\text{dom}(F))$.

As a consequence, for any $h \in \text{aff}(\text{dom}(F))$, there exists $\hat{g} \in \partial G(h)$, which therefore satisfies

$$dF(x, \tilde{h}) \geq dF(x, h) + \hat{g}^T(\tilde{h} - h)$$

for any $\tilde{h} \in \mathbb{R}^d$ (the upper bound is infinite if $\tilde{h} \notin \text{dom}(G)$). Letting $\tilde{h} \to 0$, we get $dF(x, h) \leq \hat{g}^T h$.

Also, by positive homogeneity, we have

$$t\, dF(x, \tilde{h}) \geq dF(x, h) + \hat{g}^T(t\tilde{h} - h)$$

for all $t > 0$, which requires $dF(x, \tilde{h}) \geq \hat{g}^T \tilde{h}$ for all $\tilde{h}$, and in particular $dF(x, h) = \hat{g}^T h$.

Since

$$F(x + \tilde{h}) - F(x) \geq dF(x, \tilde{h}) \geq \hat{g}^T \tilde{h}$$

we see that $\hat{g} \in \partial F(x)$, with $\hat{g}^T h = dF(x, h)$, which concludes the proof.  ∎

The next proposition gives a criterion for a vector $g$ to belong to $\partial F(x)$ based on directional derivatives.

**Proposition 3.52** *Assume that $x \in \text{dom}(F)$ where $F$ is convex. If $g \in \mathbb{R}^d$ is such that*

$$dF(x, h) \geq g^T h$$

*for all $h \in \mathbb{R}^d$, then $g \in \partial F(x)$.*

PROOF  Just use the fact that $dF(x, h) \leq F(x + h) - F(x)$.  ∎

### 3.5.4  Subgradient descent

When $F$ is a non-differentiable a convex function, directions $g$ such that $-g \in \partial F(x)$ do not always provide directions of descent. Indeed, $g \in \partial F(x)$ implies

$$F(x - \alpha g) \geq F(x) - \alpha |g|^2$$

but the inequality goes in the "wrong direction." However, we know that, for any $h \in \mathbb{R}^d$, there exists $g_h \in \partial F(x)$ such that

$$dF(x, -h) = -g_h^T h \geq -g^T h$$

for all $g \in \partial F(x)$. As a consequence, any non-vanishing solution of the equation $h = g_h$ will provide a direction of descent. This suggests looking for $h \in \partial F(x)$ such that $h \neq 0$ and $|h|^2 \leq g^T h$ for all $g \in \partial F(x)$. Since $g^T h \leq |g||h|$, this requires that $|h| \leq |g|$ for all $g \in \partial F(x)$, i.e.,

$$h = \underset{\partial F(x)}{\operatorname{argmin}}(g \mapsto |g|). \tag{3.44}$$

Conversely, if $h$ is the minimal-norm element of $\partial F(x)$ (which is necessarily unique since the norm is strictly convex and $\partial F(x)$ is convex and compact), then $|h|^2 \leq |h + t(g - h)|^2$ for all $g \in \partial F(x)$ and $t \in [0, 1]$, and taking the difference yields

$$2th^T(g - h) + t^2|g - h|^2 \geq 0.$$

The fact that this holds for all $t \geq 0$ requires that $h^T(g - h) \geq 0$ as required. We have therefore proved that $h$ defined by (3.44) is a descent direction for $F$ at $x$ (it is actually the steepest descent direction: see [204] for a proof), justifying the algorithm

$$x_{t+1} = x_t - \alpha_t \underset{\partial F(x)}{\operatorname{argmin}}(g \mapsto |g|)$$

as subgradient descent iterations.

**Example.** Consider the minimization of

$$F(x) = \psi(x) + \lambda \sum_{i=1}^{n} |x^{(i)}|$$

where $\psi$ is a $C^1$ convex function on $\mathbb{R}^d$. Let $\mathcal{A}(x) = \{i : x^{(i)} = 0\}$. Then

$$\partial F(x) = \left\{ \nabla \psi(x) + \lambda \sum_{i \notin \mathcal{A}(x)} \operatorname{sign}(x^{(i)}) + \lambda \sum_{i \in \mathcal{A}(x)} \rho_i \mathfrak{e}_i : |\rho_i| \leq 1, i \in \mathcal{A}(x) \right\}$$

where $\mathfrak{e}_i$ is the $i$th vector of the canonical basis of $\mathbb{R}^d$.

For $g = \nabla \psi(x) + \lambda \sum_{i \notin \mathcal{A}(x)} \operatorname{sign}(x^{(i)}) + \lambda \sum_{i \in \mathcal{A}(x)} \rho_i \mathfrak{e}_i$, we have

$$|g|^2 = \sum_{i \notin \mathcal{A}(x)} (\partial_i F(x) + \lambda \operatorname{sign}(x^{(i)}))^2 + \sum_{i \in \mathcal{A}(x)} (\partial_i \psi(x) - \lambda \rho_i)^2.$$

Define

$$s(t) = \operatorname{sign}(t) \min(|t|, 1).$$

Then $h$ satisfying (3.44) is given by

$$h^{(i)} = \begin{cases} \partial_i \psi(x) - \lambda \operatorname{sign}(x^{(i)}) & \text{if } i \notin \mathcal{A}(x) \\ \lambda s(\partial_i \psi(x)/\lambda) & \text{if } i \in \mathcal{A}(x). \end{cases}$$

In more complex situations, the extra minimization step at each iteration of the algorithm can be challenging computationally. The following subgradient method uses an averaging approach to minimize $F$ without requiring finding subgradients with minimal norms. It simply defines

$$x_{t+1} = x_t - \alpha_t g_t, \quad g_t \in \partial F(x_t)$$

and computes

$$\bar{x}_t = \frac{\sum_{j=1}^t \alpha_j x_j}{\sum_{j=1}^t \alpha_j}.$$

We refer to [204] for a proof of convergence of this method.

### 3.5.5   Proximal Methods

**Proximal operator.**   We start with a few simple facts. Let $F$ be a closed convex function and $\psi$ be convex and differentiable, with $\mathrm{dom}(\psi) = \mathbb{R}^d$. Let $G = F + \psi$. Then $G$ is a closed convex function. Indeed, consider the sub-level set $\Lambda_a(G) = \{x : G(x) \le a\}$ and assume that $x_n \to x$ with $x_n \in \Lambda_a(g)$. Then $\psi(x_n) \to \psi(x)$ by continuity, and for all $\epsilon > 0$, we have, for large enough $n$, $F(x_n) \le a - \psi(x) + \epsilon$. This inequality remains true at the limit because $F$ is closed, yielding $G(x) \le a + \epsilon$ for all $\epsilon > 0$, so that $x \in \Lambda_a(G)$.

We have $\mathrm{ridom}(F) \cap \mathrm{ridom}(\psi) \ne \emptyset$ so that (by theorem 3.45 and proposition 3.42) $\partial G(x) = \nabla \psi(x) + \partial F(x)$. In particular, $x^*$ is a minimizer of $G$ if and only if $-\nabla \psi(x^*) \in \partial F(x^*)$.

It one assumes that $\psi$ is strongly convex, so that there exists $m$ and $L$ such that

$$\frac{m}{2}|y - x|^2 \le \psi(y) - \psi(x) - \nabla \psi(x)^T (y - x) \le \frac{L}{2}|y - x|^2$$

for all $x, y \in \mathbb{R}^d$, then a minimizer of $G$ exists and is unique. To see this, fix $x_0 \in \mathrm{ridom}(F)$ and consider the closed convex set

$$\Omega_0 = \Lambda_{G(x_0)}(G) = \{x : G(x) \le G(x_0)\}.$$

Any minimizer of $G$ must clearly belong to $\Omega_0$. If $x \in \Omega_0$, we have

$$F(x) + \psi(x_0) + \nabla \psi(x_0)^T (x - x_0) + \frac{m}{2}|x - x_0|^2 \le G(x) \le G(x_0).$$

Moreover, there exists (from proposition 3.44) an element $g \in \partial F(x_0)$ so that $F(x) \ge F(x_0) + g^T (x - x_0)$ for all $x \in \mathbb{R}^d$. We therefore get

$$F(x_0) + \psi(x_0) + (g + \nabla \psi(x_0))^T (x - x_0) + \frac{m}{2}|x - x_0|^2 \le G(x_0).$$

for all $x \in \Omega_0$, which shows that $\Omega_0$ must be bounded and therefore compact. There exists a minimizer $x^*$ of $G$ on $\Omega_0$, and therefore on all $\mathbb{R}^d$. This minimizer is unique, since the sum of a convex function and a strictly convex function is strictly convex.

In particular, for any closed convex $F$, we can apply the previous remarks to

$$G : v \mapsto F(v) + \frac{1}{2}|x - v|^2$$

where $x \in \mathbb{R}^d$ is fixed. The function $\psi : v \mapsto |v - x|^2/2$ is strongly convex (with $L = m = 1$) and $G$ therefore has a unique minimizer $v^*$. This is summarized in the following definition.

**Definition 3.53** *Let $F$ be a closed convex function. The proximal operator associated to $F$ is the mapping $\mathrm{prox}_F : \mathbb{R}^d \to \mathrm{dom}(F)$ defined by*

$$\mathrm{prox}_F(x) = \underset{\mathbb{R}^d}{\mathrm{argmin}}(v \mapsto F(v) + \frac{1}{2}|x - v|^2). \tag{3.45}$$

From the previous discussion, we also deduce

**Proposition 3.54** *Let $F$ be a closed convex function and $\alpha > 0$. We have $x' = \mathrm{prox}_{\alpha F}(x)$ if and only if $x \in x' + \alpha \partial F(x')$. In particular, $x^*$ is a minimizer of $F$ if and only if $x^* = \mathrm{prox}_{\alpha F}(x^*)$*

Let us take a few examples.

• Let $F(x) = \lambda|x|$, $x \in \mathbb{R}^d$, for some $\lambda > 0$. Then $F$ is differentiable everywhere except at $x = 0$ and $\mathrm{dom}(F) = \mathbb{R}^d$. We have $\partial F(x) = \lambda x/|x|$ for $x \neq 0$. A vector $g$ belongs to $\partial F(0)$ if and only if

$$g^T x \leq \lambda|x|$$

for all $x \in \mathbb{R}^d$, which is equivalent to $|g| \leq \lambda$ so that $\partial F(0) = \bar{B}(0, \lambda)$.

We have $x' = \mathrm{prox}_F(x)$ if and only if $x' \neq 0$ and $x = x' + \lambda x'/|x'|$ or $x' = 0$ and $|x| \leq \lambda$. For $|x| > \lambda$, the equation $x = x' + \lambda x'/|x'|$ is solved by

$$x' = \frac{|x| - \lambda}{|x|}x$$

yielding

$$\mathrm{prox}_F(x) = \begin{cases} \dfrac{|x| - \lambda}{|x|}x \text{ if } |x| \geq \lambda \\ 0 \text{ otherwise} \end{cases} \tag{3.46}$$

• Let $\Omega$ be a closed convex set. Then $\mathrm{prox}_{\sigma_\Omega} = \mathrm{proj}_\Omega$, the projection operator on $\Omega$, as directly deduced from the definition.

The following proposition can then be compared to proposition 3.38.

**Proposition 3.55** *Let F be a closed convex function. Then* $\text{prox}_F$ *is 1-Lipschitz: for all* $x, y \in \mathbb{R}^d$,

$$|\text{prox}_F(x) - \text{prox}_F(y)| \leq |x - y|. \qquad (3.47)$$

PROOF Let $x' = \text{prox}_F(x)$ and $y' = \text{prox}_F(y)$. Then, there exists $g \in \partial F(x')$ and $h \in \partial_F(y')$ such that $x = x' + g$ and $y = y' + h$. Moreover, we have

$$F(y') - F(x') \geq g^T(y' - x')$$
$$F(x') - F(y') \geq h^T(x' - y')$$

from which we deduce $g^T(y' - x') \leq h^T(y' - x')$ or $(h - g)^T(x' - y') \geq 0$. Expressing $g, h$ in terms or $x, x', y, y'$, we get $(y - x - y' + x')^T(y' - x') \geq 0$ or

$$|y' - x'|^2 \leq (y - x)^T(y' - x') \leq |y - x||y' - x'|$$

which is only possible if $|y' - x'| \leq [y - x|$. ∎

If $F$ is differentiable, then $x' = \text{prox}_{\alpha F}(x)$ satisfies

$$x' = x - \alpha \nabla F(x')$$

so that $x \mapsto \text{prox}_{\alpha F}(x)$ can be interpreted as an implicit version of the standard gradient step $x \mapsto x - \alpha \nabla F(x)$. The iterations $x(t+1) = \text{prox}_{\alpha_t F}(x(t))$ provide an algorithm that converges to a minimizer of $F$ (this will be justified below). This algorithm is rarely practical, however, since the minimization required at each step is not necessarily much easier to perform than minimizing $F$ itself. The proximal operator, however, is especially useful when combined with splitting methods.

**Proximal gradient descent.** Assume that the objective function $F$ takes the form

$$F(x) = G(x) + H(x) \qquad (3.48)$$

where $G$ is $C^1$ on $\mathbb{R}^d$ and $H$ is a closed convex function. We first note that

$$dF(x, h) = \lim_{t \downarrow 0} \frac{F(x + th) - F(x)}{t}$$

is well defined (even if $G$ is not convex, because it is smooth), with

$$dF(x, h) = \nabla G(x)^T h + dH(x, h)$$

In particular, if $x^*$ be a minimizer of $F$, then $dF(x,h) \geq 0$ for all $h$ so that $dH(x,h) \geq -\nabla G(x)^T h$ for all $h$. Using proposition 3.52, this shows that $-\nabla G(x) \in \partial H(x)$, which is a necessary condition for optimality for $F$ (which is sufficient if $G$ is convex).

Proximal gradient descent implements the algorithm

$$x_{t+1} = \text{prox}_{\alpha_t H}(x_t - \alpha_t \nabla G(x_t)). \tag{3.49}$$

We note that a stationary point of this algorithm, i.e. a point $x$ such that $x = \text{prox}_{\alpha_t H}(x - \alpha_t \nabla G(x))$ must be such that $x - \alpha_t \nabla G(x) \in x + \alpha_t \partial H(x)$, so that $-\nabla G(x) \in \partial H(x)$. This shows that the property of being stationary does not depend on $\alpha_t > 0$, and is equivalent to the necessary optimality condition that was just discussed.

We first study this algorithm under the assumption that $G$ is $L$-$C^1$, which implies that, for all $x, y \in \mathbb{R}^d$.

$$G(y) \leq G(x) + \nabla G(x)^T (y - x) + \frac{L}{2}|x - y|^2.$$

At iteration $t$, we have

$$x_t - \alpha_t \nabla G(x_t) \in x_{t+1} + \alpha_t \partial H(x_{t+1})$$

which implies, in particular

$$\alpha_t H(x_t) - \alpha_t H(x_{t+1}) \geq (x_t - x_{t+1})^T (x_t - x_{t+1} - \alpha_t \nabla G(x_t)) = |x_t - x_{t+1}|^2 + \alpha_t \nabla G(x_t)^T (x_{t+1} - x_t)$$

Dividing by $\alpha_t$ and adding $G(x_t) - G(x_{t+1})$, we get

$$F(x_t) - F(x_{t+1}) \geq \frac{1}{\alpha_t}|x_t - x_{t+1}|^2 + G(x_t) + \nabla G(x_t)^T (x_{t+1} - x_t) - G(x_{t+1})$$

$$\geq \left(\frac{1}{\alpha_t} - \frac{L}{2}\right)|x_t - x_{t+1}|^2 \tag{3.50}$$

so that proximal gradient descent iterations reduce the objective function as soon as $\alpha_t \leq 2/L$.

Assuming that $\alpha_t < 2/L$, (3.50) can be rewritten as

$$\left|\frac{x_{t+1} - x_t}{\alpha_t}\right|^2 \leq \frac{2}{\alpha_t(2 - \alpha_t L)}(F(x_t) - F(x_{t+1})).$$

This inequality should be compared to (3.11) in the unconstrained case. It yields, in particular, the inequality

$$\min\left\{\left|\frac{x_{t+1} - x_t}{\alpha_t}\right| : t \leq T\right\} \leq \frac{F(x_0) - \min F}{2T \min\{\alpha_t(2 - \alpha_t L), t \leq T\}}. \tag{3.51}$$

As a consequence, if one runs proximal gradient descent until $|x_{t+1} - x_t|/\alpha_t$ is small enough, the algorithm will terminate in finite time as soon as $\alpha_t$ is bounded from below (and, in particular, if $\alpha_t$ is constant).

If we assume that $G$ is convex, in addition to being $L$-$C^1$, then we have a stronger result. Let $x^*$ be a minimizer of $F$. Then, using again $x_t - \alpha_t \nabla G(x_t) \in x_{t+1} + \alpha_t \partial H(x_{t+1})$, we have

$$\alpha_t H(x^*) - \alpha_t H(x_{t+1}) \geq (x^* - x_{t+1})^T (x_t - x_{t+1} - \alpha_t \nabla G(x_t))$$

and

$$\begin{aligned}
\alpha_t F(x^*) - \alpha_t F(x_{t+1}) &\geq (x^* - x_{t+1})^T (x_t - x_{t+1}) - \alpha_t (x^* - x_{t+1})^T \nabla G(x_t)) + \alpha_t G(x^*) - \alpha_t G(x_{t+1}) \\
&\geq (x^* - x_{t+1})^T (x_t - x_{t+1}) - \alpha_t (x^* - x_t)^T \nabla G(x_t)) + \alpha_t G(x^*) \\
&\quad + \alpha_t (x_{t+1} - x_t)^T \nabla G(x_t) - \alpha_t G(x_{t+1}) \\
&\geq (x^* - x_{t+1})^T (x_t - x_{t+1}) - \frac{\alpha_t L}{2} |x_t - x_{t+1}|^2
\end{aligned}$$

Assuming that $\alpha_t L \leq 1$, then

$$\alpha_t F(x^*) - \alpha_t F(x_{t+1}) \geq (x^* - x_{t+1})^T (x_t - x_{t+1}) - \frac{1}{2}|x_t - x_{t+1}|^2 = \frac{1}{2}(|x_{t+1} - x^*|^2 - |x_t - x^*|^2),$$

which we rewrite as

$$\alpha_t (F(x_{t+1}) - F(x^*)) \leq \frac{1}{2}(|x_t - x^*|^2 - |x_{t+1} - x^*|^2)$$

Note that, from (3.50), we also have

$$F(x_{t+1}) \leq F(x_t) - \frac{1}{2\alpha_t}|x_t - x_{t+1}|^2$$

when $\alpha_t L \leq 1$, which shows, in particular that $F(x_t)$ is decreasing. Fixing a time $T$, we have, from these two observations

$$\alpha_t (F(x_T) - F(x^*)) \leq \frac{1}{2}(|x_t - x^*|^2 - |x_{t+1} - x^*|^2)$$

for all $t \leq T - 1$, and summing over $T$,

$$(F(x_T) - F(x_*)) \sum_{t=0}^{T-1} \alpha_t \leq \frac{1}{2}(|x_0 - x^*|^2 - |x_T - x^*|^2)$$

yielding

$$F(x_T) - F(x_*) \leq \frac{|x_0 - x^*|^2}{2\sum_{t=0}^{T-1} \alpha_t}. \tag{3.52}$$

We summarize this in the following theorem, specializing to the case of constant step $\alpha_t$.

**Theorem 3.56** *Let G be am L-$C^1$ function defined on $\mathbb{R}^d$ and H be closed convex. Assume that $F = G+H$ has a minimizer $x^*$. Then the algorithm (3.49) run with $\alpha_t = \alpha \leq 1/L$ for all t is such that, for all $T > 0$,*

$$F(x_T) - F(x_*) \leq \frac{|x_0 - x^*|^2}{2\alpha T}. \tag{3.53}$$

Also, when $G = 0$, $F = H$, we retrieve the proximal iterations algorithm

$$x_{t+1} = \text{prox}_{\alpha F}(x_t), \tag{3.54}$$

and we have just proved that it converges for any $\alpha > 0$ as soon as $F$ is a closed convex function.

One gets a stronger result under the assumption that $G$ is $C^2$, and is such that the eigenvalues of $\nabla^2 G(x)$ are included in a fixed interval $[m, L]$ for all $x \in \mathbb{R}^d$ with $m > 0$. Such a $G$ is strongly convex, which implies that $F$ has a unique minimizer. We have

$$|x_{t+1} - x^*| = \left|\text{prox}_{\alpha_t H}(x_t - \alpha_t \nabla G(x_t)) - \text{prox}_{\alpha_t H}(x^* - \alpha_t \nabla G(x^*)\right|$$
$$\leq |x_t - x^* - \alpha_t(\nabla G(x_t)) - \nabla G(x^*))|.$$

Write

$$|x_t - x^* - \alpha_t(\nabla G(x_t)) - \nabla G(x^*)| = \left|\int_0^1 (\text{Id}_{\mathbb{R}^n} - \alpha_t \nabla^2 G(x^* + t(x_t - x^*)))(x_t - x^*)dt\right|$$
$$\leq \int_0^1 \left|(\text{Id}_{\mathbb{R}^n} - \alpha_t \nabla^2 G(x^* + t(x_t - x^*)))(x_t - x^*)\right| dt$$
$$\leq \max(|1 - \alpha_t m|, |1 - \alpha_t L|)|x_t - x^*|$$

where we have use the fact that the eigenvalues of $\text{Id}_{\mathbb{R}^n} - \alpha_t \nabla^2 G(x)$ are included in $[1 - \alpha_t L, 1 - \alpha_t m]$ for all $x \in \mathbb{R}^d$. If one assumes that $\alpha_t \leq 1/L$, so that $\max(|1 - \alpha_t m|, |1 - \alpha_t L|) \leq 1 - \alpha_t m$, one gets

$$|x_{t+1} - x^*| \leq (1 - \alpha_t m)|x_t - x^*|.$$

Iterating this inequality, we get the theorem that we state for constant $\alpha_t$.

**Theorem 3.57** *Let $F = G + H$ where G is a $C^2$ convex function and H is a closed convex function. Assume that the eigenvalues of $\nabla^2 G$ are uniformly included in $[m, L]$ with $m > 0$. Let $x^*$ argmin $F$.*

*Let $(x_t)$ satisfy (3.49) with constant $\alpha_t = \alpha < 1/L$. Then*

$$|x_t - x^*| \leq (1 - \alpha m)^t |x_0 - x^*|.$$

Note that these results also apply to projected gradient descent (section 3.4.4), which is a special case (taking $G = \sigma_\Omega$).

## 3.6  Duality

### 3.6.1  Generalized KKT conditions

A constrained convex minimization problem consists in the minimization of a closed convex function $F$ over a closed convex set $\Omega \subset \mathrm{ridom}(F)$. We have seen in theorem 3.36 that, for smooth $F$, any solution $x^*$ of this problem had to satisfy $-\nabla F(x^*) \in \mathcal{N}_\Omega(x)$ where

$$\mathcal{N}_\Omega(x) = \{h : h^T(y - x) \leq 0 \text{ for all } y \in \Omega\}.$$

The next theorem generalizes this property to the non-smooth convex case, for which the necessary optimality condition is also sufficient.

**Theorem 3.58** *Let $F$ be a closed convex function, $\Omega \subset \mathrm{ridom}(F)$ a nonempty closed convex set. Then $x^* \in \mathrm{argmin}_\Omega F$ if and only if*

$$0 \in \partial F(x^*) + \mathcal{N}_\Omega(x^*)$$

PROOF Introduce the indicator function $\sigma_\Omega$. Then minimizing $F$ over $\Omega$ is the same as minimizing $G = F + \sigma_\Omega$ over $\mathbb{R}^d$. The assumptions imply that $\mathrm{ridom}(\sigma_\Omega) = \mathrm{relint}(\Omega) \subset \mathrm{ridom}(F)$ and therefore

$$\partial G(x) = \partial F(x) + \partial \sigma_\Omega(x)$$

for all $x \in \Omega$. Since

$$\partial \sigma_\Omega(x) = \mathcal{N}_\Omega(x)$$

the result follows for the characterization of minimum of convex functions.  ■

In the following, we will restrict to the situation in which $F$ is finite (i.e., $\mathrm{dom}(F) = \mathbb{R}^d$) and $\Omega$ is defined through a finite number of equalities and inequalities, taking the form

$$\Omega = \left\{x \in \mathbb{R}^d : \gamma_i(x) = 0, i \in \mathcal{E} \text{ and } \gamma_i(x) \leq 0, i \in \mathcal{I}\right\}$$

for functions $(\gamma_i, i \in \mathcal{C} = \mathcal{E} \cup \mathcal{I})$ such that $\gamma_i : x \mapsto b_i^T x + \beta_t$ is affine for all $i \in \mathcal{E}$ and $\gamma_i$ is closed convex for all $i \in \mathcal{I}$. This is similar to the situation considered in section 3.4.1, with additional convexity assumptions, but without assuming smoothness. We recall the definition of active constraints from section 3.4.1, namely, for $x \in \Omega$,

$$\mathcal{A}(x) = \{i \in \mathcal{C} : \gamma_i(x) = 0\}.$$

Following the discussion in the smooth case, define the set $\mathcal{N}'_\gamma(x) \subset \mathbb{R}^d$ by

$$\mathcal{N}'_\gamma(x) = \left\{\sum_{i \in \mathcal{A}(x)} \lambda_i s_i : s_i \in \partial \gamma_i(x), i \in \mathcal{A}(x), \lambda_i \geq 0, i \in \mathcal{A}(x) \cap \mathcal{I}\right\}.$$

The property $0 \in \partial F(x^*) + \mathcal{N}_\gamma(x^*)$ is the expression of the KKT conditions in the non-smooth case. It holds for $x^* \in \operatorname{argmin}_\Omega F$ as soon as $\mathcal{N}_\Omega(x^*) = \mathcal{N}'_\gamma(x^*)$, which is true under appropriate constraint qualifications. We here replace the MF-CQ in definition 3.30 by the following conditions that do not involve gradients.

**Definition 3.59** *Let $(\gamma_i, i \in \mathcal{C} = \mathcal{E} \cup \mathcal{I})$ be a set of equality and inequality constraints, with $\gamma_i : x \mapsto b_i^T x + \beta_i$, $i, \in \mathcal{E}$ and $\gamma_i$ closed convex, $i \in \mathcal{I}$. One says that these constraints satisfy the Slater constraint qualifications (Sl-CQ) if and only if:*

*(Sl 1)  The vectors $(b_i, i \in \mathcal{E})$ are linearly independent.*

*(Sl 2)  There exists $x \in \mathbb{R}^d$ such that $\gamma_i(x) = 0$ for $i \in \mathcal{E}$ and $\gamma_i(x) < 0$ for $i \in \mathcal{I}$.*

The first constraint is a very mild condition. When it is not satisfied, this means that some $b_i$'s are linear combinations of others, and equality constraints for the latter implies equality constraints for the former. These redundancies can therefore be removed without changing the problem.

Note that (Sl2) can be replaced by the apparently weaker condition that, for all $i \in \mathcal{I}$, there exists $x_i \in \mathbb{R}^d$ satisfying all the constraints and $\gamma_i(x_i) < 0$. Indeed, if this is true, then the average, $\bar{x}$, of $(x_i, i \in \mathcal{I})$ also satisfies the equality constraints by linearity, and if $i \in \mathcal{I}$,

$$\gamma_i(\bar{x}) \leq \frac{1}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} \gamma_i(x^{(j)}) \leq \frac{1}{|\mathcal{I}|} \gamma_i(x^{(i)}) < 0.$$

The following proposition makes a connection between the Slater conditions and the MF-CQ in definition 3.30.

**Proposition 3.60** *Assume that $\gamma_i$, $i \in \mathcal{I}$ are convex $C^1$ functions. Then, if there exists a feasible point $x^*$ that satisfies the MF-CQ, there exists another point $x$ satisfying the Sl-CQ. Conversely, if there exists $x$ satisfying the Sl-CQ, then every feasible point $x^*$ satisfies the MF-CQ.*

PROOF The linear independence conditions on equality constraints are the same in MF-CQ and Sl-CQ, so that we only need to consider inequality constraints.

Let $x^*$ satisfy MF-CQ, and take $h \neq 0$ such that $b_i^T h = 0$ for all $i \in \mathcal{E}$, and $\nabla \gamma_i(x^*)^T h < 0$, $i \in \mathcal{A}(x) \cap \mathcal{I}$. Then $x^* + th$ satisfies the equality constraints for all $t \in \mathbb{R}$. If $i \in \mathcal{I}$ is not active, then $\gamma_i(x^*) < 0$ and this will remain true at $x^* + th$ for small $t$ by continuity. If $i \in \mathcal{A}(x) \cap \mathcal{I}$, then a first order expansion gives $\gamma_i(x^* + th) = t \nabla \gamma_i(x^*)^T h + o(|h|)$, which is also negative for small enough $t > 0$. So, $x^* + th$ satisfies the Sl-CQ for small enough $t > 0$.

Conversely, let $x$ satisfy the Sl-CQ. Take a feasible point $x^*$. If $x^* = x$, then there is no active inequality constraint and $x^*$ satisfies MF-CQ. Assume $x^* \neq x$ and let $h = x - x^*$. Then $b_i^T h = 0$ for all $i \in \mathcal{E}$, and if $i \in \mathcal{I} \cap \mathcal{A}(x^*)$,

$$0 > \gamma_i(x) = \gamma_i(x^* + h) \geq \gamma_i(x^*) + \nabla \gamma_i(x^*)^T h = \nabla \gamma_i(x^*)^T h$$

so that $x^*$ satisfies MF-CQ.                                                    ■

The following theorem, that we give without proof, states that the Slater conditions implies that the KKT conditions are satisfied for a minimizer.

**Theorem 3.61** *Assume that all the constraints are affine, or that they satisfy the Sl-CQ in definition 3.59. Let $x^* \in \text{argmin}_\Omega F$. Then $\mathcal{N}_\Omega(x^*) = \mathcal{N}'_\gamma(x^*)$, so that there exist $s_0 \in \partial F(x^*)$, $s_i \in \partial \gamma_i(x^*)$, $i \in \mathcal{A}(x^*)$, $(\lambda_i, i \in \mathcal{A}(x^*))$ with $\lambda_i \geq 0$ if $i \in \mathcal{I} \cap \mathcal{A}(x^*)$, such that*

$$s_0 + \sum_{i \in \mathcal{A}(x)} \lambda_i s_i = 0 \tag{3.55}$$

### 3.6.2   Dual problem

Consider the Lagrangian

$$L(x, \lambda) = F(x) + \sum_{i \in \mathcal{C}} \lambda_i \gamma_i(x)$$

defined in (3.29) and let $D = \{\lambda : \lambda_i \geq 0, i \in \mathcal{I}\}$. Because the functions $\gamma_i$ are nonpositive on $\Omega$, we have

$$L(x, \lambda) \leq F(x)$$

for all $x \in \Omega$ and $\lambda \in D$, which implies that

$$L^*(\lambda) = \inf\{L(x, \lambda) : x \in \mathbb{R}^d\}$$

is such that $L^*(\lambda) \leq F(x)$ for all $\lambda \in D$ and $x \in \Omega$. Define

$$\hat{d} = \sup\{L^*(\lambda) : \lambda \in D\}$$

and

$$\hat{p} = \inf\{F(x) : x \in \Omega\},$$

whose computations respectively represent the *dual* and *primal* problems. Then, we have $\hat{d} \leq \hat{p}$.

We did not need much of our assumptions (not even $F$ to be convex) to reach this conclusion. When the converse inequality is true (so that the *duality gap* $\hat{p} - \hat{d}$ vanishes), the dual problem provides important insights on the primal problem, as well as alternative ways to solve it. This is true under the Slater conditions.

**Theorem 3.62** *The duality gap vanishes when the constraints are all affine, or when they satisfy the Sl-CQ in definition 3.59. In this case, any solution $\lambda^*$ of the dual problem provides Lagrange multipliers in theorem 3.61 and conversely.*

We justify this statement, as a consequence of theorem 3.61 and the following analysis. The Lagrangian $L(x, \lambda)$ is linear in $\lambda$, and when $\lambda \in D$, is a convex function of $x$. Moreover, one can use subdifferential calculus (theorem 3.45) to conclude that, for any $\lambda \in D$, (3.55) expresses the fact that $0 \in \partial_x L(x^*, \lambda)$, i.e., that $x^* \in \operatorname{argmin}_{\mathbb{R}^d} L(\cdot, \lambda)$.

Fixing $x \in \mathbb{R}^d$, one can also consider the maximization of $L$ in $\lambda \in D$. Clearly, if $x \notin \Omega$, so that $\gamma_i(x) \neq 0$ for some $i \in \mathcal{E}$ or $\gamma_i(x) > 0$ for some $i \in \mathcal{I}$, then $\max_D L(x, \lambda) = +\infty$. If $x \in \Omega$, then the slackness conditions, which require $\lambda^{(i)} \gamma_i(x) = 0$, $i \in \mathcal{I}$, ensure that $\lambda \in \operatorname{argmax}_D L(x, \cdot)$.

As a consequence, any pair $x^* \in \Omega$, $\lambda^* \in D$ satisfying the KKT conditions is such that

$$L(x^*, \lambda) \leq L(x^*, \lambda^*) \leq L(x, \lambda^*) \tag{3.56}$$

for all $x \in \mathbb{R}^d$ and $\lambda \in D$. Such a pair $(x^*, \lambda^*)$ is called a *saddle point* of the function $L$. Conversely, any saddle point of $L$, i.e., any $(x^*, \lambda^*) \in \mathbb{R}^d \times D$ satisfying (3.56), must be such that $x^* \in \Omega$ (to ensure that $L(x^*, \cdot)$ is bounded), and satisfies the KKT conditions.

We therefore obtain the equivalence of the two properties, for $(x^*, \lambda^*) \in \mathbb{R}^d \times D$:

(i) $x^* \in \Omega$ and $(x^*, \lambda^*)$ satisfies the KKT conditions.

(ii) Equation (3.56) holds for all $(x, \lambda) \in \mathbb{R}^d \times D$.

Consider now the additional condition that

(iii) $x^* \in \operatorname{argmin}_\Omega F$ and $\lambda^* \in \operatorname{argmax}_D L^*$.

We already know that, if $(x^*, \lambda^*)$ satisfy the KKT conditions, then $x^* \in \operatorname{argmin}_\Omega F$ (because $\mathcal{N}'_\gamma(x^*) \subset \mathcal{N}_\Omega(x^*)$). Moreover, if (3.56) holds, then the inequality $L(x^*, \lambda) \leq L(x^*, \lambda^*)$ implies that $L^*(\lambda) \leq L(x^*, \lambda^*)$ for all $\lambda \in D$. The inequality $L(x^*, \lambda^*) \leq L(x, \lambda^*)$ for all $x$ implies that $L(x^*, \lambda^*) \leq L^*(\lambda^*)$. We therefore obtain the fact that $\lambda^* \in \operatorname{argmax} L^*(\lambda)$. To summarize, we have

$$\text{(i)} \Leftrightarrow \text{(ii)} \Rightarrow \text{(iii)}.$$

To obtain the final equivalence, we need to assume constraints qualifications, such as Slater's conditions, to ensure that $\mathcal{N}'_\gamma(x^*) = \mathcal{N}_\Omega(x^*)$. If this holds, then (iii) implies (via theorem 3.61) that there exists $\tilde{\lambda}$ such that (i) and (ii) are satisfied for

$(x^*, \tilde{\lambda})$, with $L(x^*, \tilde{\lambda}) = L^*(\tilde{\lambda})$ and $\tilde{\lambda} \in \text{argmin}_D L^*$. This shows that $L^*(\tilde{\lambda}) = L^*(\lambda^*)$. Moreover, from (3.56), we have

$$L(x^*, \lambda^*) \leq L(x^*, \tilde{\lambda}) = L^*(\tilde{\lambda}),$$

and, by definition of $L^*$, $L(x^*, \lambda^*) \geq L^*(\lambda^*)$. This shows that $L(x^*, \lambda^*) = L(x^*, \tilde{\lambda})$. As a consequence, for all $(x, \lambda) \in \mathbb{R}^d \times D$:

$$L(x^*, \lambda) \leq L(x^*, \tilde{\lambda}) = L(x^*, \lambda^*) = L^*(\lambda^*) = \inf_{\mathbb{R}^d} L(\cdot, \lambda^*) \leq L(x, \lambda^*)$$

so that $(x^*, \lambda^*)$ satisfies (ii).

### 3.6.3   Example: Quadratic programming

Quadratic programming problems minimize $F(x) = \frac{1}{2} x^T A x - b^T x$, where $A$ is a positive semidefinite matrix and $b \in \mathbb{R}^d$, subject to affine constraints $c_i^T x - d_i = 0$, $i \in \mathcal{E}$ and $c_i^T x - d_i \leq 0$, $i \in \mathcal{I}$.

We here consider the following objective function. Introduce variables $x \in \mathbb{R}^d$, $x_0 \in \mathbb{R}$ and $\xi \in \mathbb{R}^N$ and minimize, for a fixed parameter $\gamma$,

$$F(x, x_0, \xi) = \frac{1}{2} |x|^2 + \gamma \sum_{k=1}^{N} \xi^{(k)}$$

subject to constraints, for $k = 1, \ldots, N$ $\xi^{(k)} \geq 0$ and

$$b_k(x_0 + x^T a_k) + \xi^{(k)} \geq 1$$

where $b_k \in \{-1, 1\}$ and $a_k \in \mathbb{R}^n$ respectively denote the $k$th output and input training sample. This algorithm minimizes a quadratic function of the input variables $(x, x_0, \xi)$ subject to linear constraints, and is an instance of a quadratic programming problem (this is actually the support vector machine problem for classification, which will be described in section 8.4.1).

Introduce Lagrange multipliers $\eta_k$ for the constraint $\xi^{(k)} \geq 0$ and $\alpha_k$ for $b_k(x_0 + x^T a_k) + \xi^{(k)} \geq 1$. The Lagrangian then takes the form

$$L(x, x_0, \xi, \alpha, \eta) = \frac{1}{2} |x|^2 + \gamma \sum_{k=1}^{N} \xi^{(k)} - \sum_{i=1}^{N} \eta_k \xi^{(k)} - \sum_{k=1}^{N} \alpha_k (b_k(x_0 + x^T a_k) + \xi^{(k)} - 1)$$

$$= \frac{1}{2} |x|^2 + \sum_{k=1}^{N} (\gamma - \eta_k - \alpha_k) \xi^{(k)} - x_0 \sum_{k=1}^{N} \alpha_k b_k - x^T \sum_{k=1}^{N} \alpha_k b_k a_k + \sum_{k=1}^{N} \alpha_k.$$

We compute the dual Lagrangian $L^*$ by minimizing with respect to the primal variables. We note that $L^*(\alpha, \eta) = -\infty$ when $\sum_{k=1}^{N}$ $alpha_k b_k \neq 0$, so that $\sum_{k=1}^{N} \alpha_k b_k = 0$ is a constraint for the dual problem. The minimization in $\xi^{(k)}$ also gives $-\infty$ unless $\gamma - \eta_k - \alpha_k = 0$, which is therefore another constraint. Finally, the optimal values of $x$ is

$$x = \sum_{k=1}^{N} \alpha_k b_k a_k$$

and we obtain the expression of the dual problem, which maximizes

$$-\frac{1}{2} \sum_{k,l=1}^{N} \alpha_k \alpha_l b_k b_l a_k^T a_l + \sum_{k=1}^{N} \alpha_k$$

subject to $\eta_k, \alpha_k \geq 0$, $\gamma - \eta_k - \alpha_k = 0$ and $\sum_{k=1}^{N} \alpha_k b_k = 0$. The conditions on $\eta_k$ and $\alpha_k$ can be rewritten as $0 \leq \alpha_k \leq \gamma$, $\eta_k = \gamma - \alpha_k$, and since the rest of the problem does not depends on $\eta$, the dual problem can be reduced to maximizing

$$L^*(\alpha) = -\frac{1}{2} \sum_{k,l=1}^{N} \alpha_k \alpha_l a_k^T a_l + \sum_{k=1}^{N} \alpha_k$$

subject to $0 \leq \alpha_k \leq \gamma$ and $\sum_{k=1}^{N} \alpha_k b_k = 0$.

### 3.6.4  Proximal iterations and augmented Lagrangian

The concave function $L^*$ can be maximized by minimizing $-L^*$ using proximal iterations ((3.54)):

$$\lambda(t+1) = \text{prox}_{-\alpha_t L^*}(\lambda(t)) = \underset{D}{\text{argmax}}(\lambda \mapsto L^*(\lambda) - \frac{1}{2\alpha_t}|\lambda - \lambda(t)|^2).$$

Introduce the function

$$\varphi(x, \lambda) = F(x) + \sum_{i \in C} \lambda^{(i)} \gamma_i(x) - \frac{1}{2\alpha_t}|\lambda - \lambda(t)|^2$$

so that

$$\lambda(t+1) = \underset{\mu \in D}{\text{argmax}} \underset{x \in \mathbb{R}^n}{\inf} \varphi(x, \mu).$$

The function $\varphi$ is convex in $x$ and strongly concave in $\mu$. Results in "minimax theory" [27] implies that one has the equality

$$\underset{\mu \in D}{\max} \underset{x \in \mathbb{R}^n}{\inf} \varphi(x, \mu) = \underset{x \in \mathbb{R}^d}{\inf} \underset{\mu \in D}{\sup} \varphi(x, \mu) \tag{3.57}$$

(Note that the left-hand side of this equation is never larger than the right-hand side, but their equality requires additional hypotheses—which are satisfied in our context—in order to hold.)

Importantly, the maximization in $\mu$ in the right-hand side has a closed form solution. It requires to maximize

$$\sum_{i \in \mathcal{C}} \left( \mu_i \gamma_i(x) - \frac{1}{2\alpha_t}(\mu_i - \lambda_i(t))^2 \right)$$

subject to $\mu_i \geq 0$ for $i \in \mathcal{I}$, and each $\mu_i$ can be computed separately. For $i \in \mathcal{E}$, there is no constraint on $\mu_i$, and one finds

$$\mu_i = \lambda_i(t) + \alpha_t \gamma_i(x),$$

and

$$\mu_i \gamma_i(x) - \frac{1}{2\alpha_t}(\mu_i - \lambda_i(t))^2 = \lambda_i(t)\gamma_i + \frac{\alpha_t}{2}\gamma_i(x)^2 = \frac{1}{2\alpha_t}(\lambda_i(t) + \alpha_t \gamma_i(x))^2 - \frac{\lambda_i(t)^2}{2\alpha_t}.$$

For $i \in \mathcal{I}$, the solution is

$$\mu_i = \max(0, \lambda_i(t) + \alpha_t \gamma_i(x))$$

and one can check that, in this case:

$$\mu_i \gamma_i(x) - \frac{1}{2\alpha_t}(\mu_i - \lambda_i(t))^2 = \frac{1}{2\alpha_t}\max(0, \lambda_i(t) + \alpha_t \gamma_i(x))^2 - \frac{\lambda_i(t)^2}{2\alpha_t}$$

As a consequence, the right-hand side of (3.57) requires to minimize

$$G(x) = F(x) + \frac{1}{2\alpha_t}\sum_{i \in \mathcal{E}}(\lambda_i(t) + \alpha_t \gamma_i(x))^2 + \frac{1}{2\alpha_t}\sum_{i \in \mathcal{I}}\max(0, \lambda_i(t) + \alpha_t \gamma_i(x)))^2$$

$$- \frac{1}{2\alpha_t}\sum_{i \in \mathcal{C}}\lambda_i(t)^2.$$

If we assume that the sub-level sets $\{x \in \Omega : F(x) \leq \rho\}$ are bounded (or empty) for any $\rho \in \mathbb{R}$, then so are the sets $\{x \in \mathbb{R}^n : G(x) \leq \rho\}$, and this is a sufficient condition for the existence of a *saddle point* for $\varphi$, which is a pair $(x^*, \lambda^*)$ such that, for all $(x, \lambda) \in \mathbb{R}^n \times D$,

$$\varphi(x^*, \lambda) \leq \varphi(x^*, \lambda^*) \leq \varphi(x, \lambda^*).$$

One can then check that this implies that $x^* \in \operatorname{argmin}_{\mathbb{R}^n} G$ while $\lambda^* = \lambda(t+1)$, so that

the latter can be computed as follows:

$$
\begin{cases}
x(t) = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ F(x) + \dfrac{1}{2\alpha_t} \sum_{i \in \mathcal{E}} (\lambda_i(t) + \alpha_t \gamma_i(x))^2 \right. \\
\qquad \left. + \dfrac{1}{2\alpha_t} \sum_{i \in \mathcal{I}} \max(0, \lambda_i(t) + \alpha_t \gamma_i(x)))^2 \right\} \\
\lambda_i(t+1) = \lambda_i(t) + \alpha_t \gamma_i(x(t)), \ i \in \mathcal{E} \\
\lambda_i(t+1) = \max(0, \lambda_i(t) + \alpha_t \gamma_i(x(t))), \ i \in \mathcal{I}
\end{cases}
\tag{3.58}
$$

These iterations define the augmented Lagrangian algorithm. Starting this algorithm with some $\lambda(0) \in \mathbb{R}^{|\mathcal{C}|}$, and constant $\alpha$, $\lambda(t)$ will converge to a solution $\hat{\lambda}$ of the dual problem. The last two iterations stabilizing imply that $\gamma_i(x(t))$ converges to $0$ for $i \in \mathcal{E}$, and also for $i \in \mathcal{I}$ such that $\hat{\lambda}_i > 0$, and that $\limsup \gamma_i(x(t)) = 0$ otherwise. This shows that, if $x(t)$ converges to a limit $\tilde{x}$, then $G(\tilde{x}) = F(\tilde{x})$. However, for any $x \in \Omega$, we have

$$
G(x(t)) \leq G(x) \leq F(x)
$$

(the proof being left to the reader), showing that $\tilde{x} \in \operatorname{argmin}_\Omega F$.

Note that the augmented Lagrangian method can also be used in non-convex optimization problems [146], requiring in that case that $\alpha$ is small enough.

### 3.6.5 Alternative direction method of multipliers

We return to a situation considered in section 3.5.5 where the function to minimize takes the form $F(x) = G(x) + H(x)$. Here, we do not assume that $G$ or $H$ is smooth, but we will need their respective proximal operators to be easy to compute.

The problem can be reformulated as a minimization with equality constraints, namely that of minimizing $\tilde{F}(x, z) = G(x) + F(z)$ subject to $x = z$. We will actually consider a more general situation, namely the problem minimizing a function $\tilde{F}(x, z)$ subject to constraints $Ax + Bz = c$ where $A$ and $B$ are respectively $d \times n$ and $d \times m$ matrices, $x \in \mathbb{R}^n$, $z \in m\mathbb{R}^m$, $c \in \mathbb{R}^d$. The augmented Lagrangian algorithm applied to this problem leads to iterate (with only equality constraints)

$$
\begin{cases}
x_t, z_t = \operatorname{argmin}\{ G(x) + F(z) + \dfrac{1}{2\alpha_t} |\lambda_t + \alpha_t (Ax + Bz - c)|^2 \} \\
\lambda_{t+1} = \lambda_t + \alpha_t (Ax_t + Bz_t - c)
\end{cases}
$$

with $\lambda_t \in \mathbb{R}^d$.

One can now consider splitting the first step in two and iterate:

$$\begin{cases} x_t = \text{argmin}\{G(x) + F(z_{t-1}) + \dfrac{1}{2\alpha_t}|\lambda_t + \alpha_t(Ax + Bz_{t-1} - c)|^2\} \\[2ex] z_t = \text{argmin}\{G(x_t) + F(z) + \dfrac{1}{2\alpha_t}|\lambda_t + \alpha_t(Ax_t + Bz - c)|^2\} \\[2ex] \lambda_{t+1} = \lambda_t + \alpha_t(Ax_t + Bz_t - c) \end{cases} \tag{3.59}$$

These iterations constitute the "alternative direction method of multipliers," or ADMM (it is also sometimes called Douglas-Rachford splitting). It is not equivalent to the augmented Lagrangian algorithm (one would need to iterate a large number of times over the first two steps before applying the third one for this), but still satisfies good convergence properties. The reader can refer to Boyd et al. [40] for a relatively elementary proof that shows that this algorithm converges as soon as, in addition to the hypotheses that were already made, the Lagrangian

$$L(x, z, \lambda) = G(x) + H(z) + \lambda^T(Ax + Bz - c)$$

has a saddle point: there exists $x^*, z^*, y^*$ such that

$$\max_y L(x^*, z^*, \lambda) = L(x^*, z^*, \lambda^*) = \min_{x,z} L(x, z, \lambda^*).$$

## 3.7   Convex separation theorems and additional proofs

We conclude this chapter by completing some of the proofs left aside when discussion convex functions. These proofs use convex separation theorems, stated below (without proof).

**Theorem 3.63 (c.f., Rockafellar [167])** *Let $\Omega_1$ and $\Omega_2$ be two nonempty convex sets with* $\text{relint}(\Omega_1) \cap \text{relint}(\Omega_2) = \emptyset$. *Then there exists $b \in \mathbb{R}^d$ and $\beta \in \mathbb{R}$ such that $b \neq 0$,* $b^T x \leq \beta$ *for all $x \in \Omega_1$ and $b^T x \geq \beta$ for all $x \in \Omega_2$, with a strict inequality for at least one* $x \in \Omega_1 \cup \Omega_2$.

**Theorem 3.64** *Let $\Omega_1$ and $\Omega_2$ be two nonempty convex sets with $\Omega_1 \cap \Omega_2 = \emptyset$ and $\Omega_1$ compact. Then there exists $b \in \mathbb{R}^n$, $\beta \in \mathbb{R}$ and $\epsilon < 0$ such that $b^T x \leq \beta - \epsilon$ for all $x \in \Omega_1$ and $b^T x \geq \beta + \epsilon$ for all $x \in \Omega_2$.*

### 3.7.1   Proof of proposition 3.44

We start with a few general remarks. If $x \in \mathbb{R}^d$, the set $\{x\}$ is convex and $\text{relint}(\{x\}) = \{x\}$. If $\Omega$ is any convex set such that $x \notin \text{relint}(\Omega)$, then theorem 3.63 implies that there exist $b \in \mathbb{R}^d$ and $\beta \in \mathbb{R}$ such that $b^T y \geq \beta \geq b^T x$ for all $y \in \Omega$ (with $b^T y > b^T x$ for

at least one $y$). If $x$ is in $\Omega \setminus (\text{relint}(\Omega))$ (so that $x$ is a point on the relative boundary of $\Omega$), then, necessarily $b^T x = \beta$ and we can write

$$b^T y \geq b^T x$$

for all $y \in \Omega$ with a strict inequality for some $y \in \Omega$. One says that $b$ and $\beta$ provide a *supporting hyperplane* for $\Omega$ at $x$.

Now, if $F$ is a convex function, with

$$\text{epi}(F) = \{(y, a) \in \mathbb{R}^d \times \mathbb{R} : F(y) \leq a\}$$

then

$$\text{relint}(\text{epi}(F)) = \{(y, a) \in \text{ridom}(F) \times \mathbb{R} : F(y) < a\}$$

(this simple fact is proved in lemma 3.65 below). In particular, if $x \in \text{dom}(F)$, then $(x, F(x))$ must be in the relative boundary of $\text{epi}(F)$. This implies that there exists $(b, b_0) \neq (0, 0) \in \mathbb{R}^d \times \mathbb{R}$ such that, for all $(y, a) \in \text{epi}(F)$:

$$b^T y + b_0 a \geq b^T x + b_0 F(x).$$

If one assumes that $x \in \text{ridom}(F)$, then, necessarily, $b_0 \neq 0$. To show this, assume otherwise, so that $b^T y \geq b^T x$ for all $y \in \text{dom}(F)$, with $b \neq 0$. We get a contradiction using the fact that, for some $\epsilon > 0$, $[y, x - \epsilon(y - x)]$ belongs to $\text{dom}(\Omega)$, because $b^T(y - x)$ cannot have a constant sign on this segment.

So $b_0 \neq 0$ and necessarily $b_0 > 0$ to ensure that $b^T y + b_0 a$ is bounded from below for all $a \geq F(y)$. Without loss of generality, we can assume $b_0 = 1$ and we get, for all $y \in \text{dom}(F)$

$$F(y) + b^T y \geq F(x) + b^T x$$

which shows that $-b \in \partial F(x)$, justifying the fact that $\partial F(x) \neq \emptyset$ for $x \in \text{ridom}(F)$.

We now state and prove the result announced above on the relative interior of the epigraph of a convex function.

**Lemma 3.65** *Let F be a convex function with epigraph*

$$\text{epi}(F) = \{(y, a) : y \in \text{dom}(F), F(y) \leq a\}.$$

*Then*

$$\text{relint}(\text{epi}(G)) = \{(y, a) : y \in \text{ridom}(F), F(y) < a\}.$$

PROOF Let $\Gamma = \{(y, a) : y \in \text{ridom}(F), F(y) < a\}$. Assume that $(y, a) \in \text{relint}(\text{epi}(F))$. Then $(y, b) \in \text{epi}(F)$ for all $b > a$ and there exists $\epsilon > 0$ such that $(y, a) - \epsilon((y, b) - (y, a)) \in \text{epi}(F)$ which requires that $F(y) \leq a - \epsilon(b - 1) < a$. Now, take $x \in \text{dom}(F)$.

Then, $(x, F(x)) \in \text{epi}(\text{dom}(F))$ and $(y, a) - \epsilon((x, F(x)) - (y, a)) \in \text{epi}(F)$ for small enough $\epsilon$, showing that $F(y - \epsilon(x - y)) \leq (1 + \epsilon)a - \epsilon F(x)$ and $y - \epsilon(x - y) \in \text{dom}(F)$. This proves that $y \in \text{ridom}(F)$ and the fact that $\text{relint}(\text{epi}(F)) \subset \Gamma$.

Take $(y, a) \in \Gamma$, and $(x, b) \in \text{epi}(F)$. We need to show that $(y - \epsilon(x - y), a - \epsilon(b - a)) \in \text{epi}(F)$ for small enough $\epsilon$, i.e., that

$$F(y - \epsilon(x - y)) \leq a - \epsilon(b - a)$$

for small enough $\epsilon$. But this is an immediate consequence of the facts that $F$ is continuous at $y \in \text{ridom}(G)$ and $F(y) < a$.                                  ■

### 3.7.2   Proof of theorem 3.45

Assume that there exists $\bar{x} \in \text{ridom}(F_1) \cap \text{ridom}(F_2)$. Take $x \in \text{dom}(F_1) \cap \text{dom}(F_2)$ and $g \in \partial(F_1 + F_2)(x)$. We want to show that $g = g_1 + g_2$ with $g_1 \in \partial F_1(x)$ and $g_2 \in \partial F_2(x)$.

By definition, we have

$$F_1(y) + F_2(y) \geq F_1(x) + F_2(x) + g^T(y - x)$$

for all $y$. We want to decompose $g$ as $g = g_1 + g_2$ with $g_1 \in \partial F_1(x)$ and $g_2 \in \partial F_2(x)$. Equivalently, we want to find $g_2 \in \mathbb{R}^d$ such that, for all $y \in \mathbb{R}^d$,

$$F_1(y) \geq F_1(x) + (g - g_2)^T(y - x)$$
$$F_2(y) \geq F_2(x) + g_2^T(y - x)$$

First note that we can replace $F_1$ by $y \mapsto F_1(y) - F_1(x) - g^T(y - x)$ and $F_2$ by $y \mapsto F_2(y) - F_2(x)$ and assume with loss of generality that $F_1(x) = F_2(x) = 0$ and $g = 0$. Making this assumption, we need to find $g_2$ such that

$$F_1(y) \geq -g_2^T(y - x)$$
$$F_2(y) \geq g_2^T(y - x)$$

for all $y \in \mathbb{R}^d$ and some $g_2 \in \mathbb{R}^d$, under the assumption that $F_1(y) + F_2(y) \geq 0$ for all $y$. Introduce the two convex sets in $\mathbb{R}^d \times \mathbb{R}$

$$\Omega_1 = \text{epi}(F_1) = \{(y, a) \in \mathbb{R}^d \times \mathbb{R} : F_1(y) \leq a\}$$
$$\Omega_2 = \{(y, a) \in \mathbb{R}^d \times \mathbb{R} : F_2(y) \leq -a\}.$$

The set $\Omega_2$ is the image of $\text{epi}(F_2)$ by the transformation $(y, a) \mapsto (y, -a)$. We have

$$\text{relint}(\Omega_1) = \text{epi}(F_1) = \{(y, a) \in \text{ridom}(F_1) \times \mathbb{R} : F_1(y) < a\}$$
$$\text{relint}(\Omega_2) = \{(y, a) \in \text{ridom}(F_2) \times \mathbb{R} : F_2(y) < -a\}.$$

Since $F_1 + F_2 \geq 0$, $\Omega_1$ and $\Omega_2$ have non-intersecting relative interiors. We can apply the first separation theorem, providing $\bar{b} = (b, b_0) \in \mathbb{R}^d \times \mathbb{R}$ and $\beta \in \mathbb{R}$ such that $\bar{b} \neq (0,0)$, $b^T y + b_0 a - \beta \leq 0$ for $(y, a) \in \Omega_1$ and $b^T y + b_0 a - \beta \geq 0$ for $(y, b) \in \Omega_2$, with a strict inequality for at least one point in $\Omega_1 \cup \Omega_2$. We therefore obtain the fact that, for all $y$ and $a$,

$$F_1(y) \leq a \Longrightarrow b^T y + b_0 a - \beta \leq 0$$
$$F_2(y) \leq -a \Longrightarrow b^T y + b_0 a - \beta \geq 0.$$

We claim that $b_0 \neq 0$. Indeed, if $b_0 = 0$, the statement for $F_1$ would imply that $b^T y - \beta \leq 0$ for all $y \in \text{dom}(F_1)$ and the one on $F_2$ that $b^T y - \beta \geq 0$ for $y \in \text{dom}(F_2)$. The point $\bar{x} \in \text{relint}(\Omega_1) \cap \text{relint}(\Omega_2)$ should then satisfy $b^T \bar{x} - \beta = 0$. We know that there exists a point $y \in \Omega_1 \cup \Omega_2$ such that $b^T y \neq \beta$. Assume that $y \in \Omega_1$, so that $b^T y - \beta < 0$ and take $\epsilon > 0$ such that $\tilde{y} = \bar{x} - \epsilon(y - \bar{x}) \in \Omega_1$. Then

$$b^T \tilde{y} - \beta = -\epsilon(b^T y - \beta) < 0,$$

which is a contradiction. A similar contradiction is obtained when $y$ belongs to $\Omega_2$, yielding the fact that $b_0$ cannot vanish.

Moreover, we clearly need $b_0 < 0$ to ensure that $b^T y + b_0 a - \beta \leq 0$ for all large enough $a$ if $y \in \text{dom}(\Omega_1)$. There is then no loss of generality in assuming $b_0 = -1$ and we get

$$F_1(y) \leq a \Longrightarrow b^T y - \beta \leq a$$
$$F_2(y) \leq -a \Longrightarrow b^T y - \beta \geq a,$$

which is equivalent to

$$-F_2(y) \leq b^T y - \beta \leq F_1(y)$$

Taking $y = x$ gives $\beta = b^T x$ and we get the desired inequality with $g_2 = -b$.

### 3.7.3   Proof of theorem 3.46

Let $\bar{x} \in \mathbb{R}^m$ such that $A\bar{x} \in \text{ridom}(F)$. We need to prove that $\partial G(x) \subset A^T \partial F(Ax + b)$ when $G(x) = F(Ax + b)$. We assume in the following that $b = 0$, since the theorem with $G(x) = F(x + b)$ is obvious. If $g \in \partial G(x)$, we have

$$F(Ay) \geq F(Ax) + g^T(y - x)$$

for all $y \in \mathbb{R}^m$. We want to show that there exists $h \in \mathbb{R}^d$ such that $g = A^T h$ and, for all $z \in \mathbb{R}^d$,

$$F(z) \geq F(Ax) + h^T(z - Ax) = F(Ax) + h^T z - g^T x.$$

Let $\Omega_1 = \text{epi}(F) = \{(z,a) :, z \in \mathbb{R}^d, F(z) \le a\}$ and

$$\Omega_2 = \{(Ay,a) : y \in \mathbb{R}^m, a = g^T(y-x) + G(x)\} \subset \mathbb{R}^d \times \mathbb{R}.$$

Note that $\Omega_2$ is an affine space with $\text{relint}(\Omega_2) = \Omega_2$. If $(z,a) \in \text{relint}(\Omega_1) \cap \Omega_2$, then $z = Ay$ for some $y \in \mathbb{R}^m$ and $g^T(y-x) + G(x) > F(z) = G(y)$. This contradicts the fact that $g \in \partial G(x)$ and shows that $\text{relint}(\Omega_1) \cap \Omega_2 = \emptyset$. As a consequence, there exist $(b, b_0) \ne (0,0)$ and $\beta$ such that

$$F(z) \le a \Rightarrow b^T z + b_0 a \le \beta$$
$$z = Ay, a = g^T(y-x) + G(x) \Rightarrow b^T z + b_0 a \ge \beta$$

Assume, to get a contradiction, that $b_0 = 0$ (so that $b \ne 0$). Then $b^T Ay \ge \beta$ for all $y$, which is only possible if $b$ is perpendicular to the range of $A$ and $\beta \le 0$. On the other hand, $F(A\bar{x}) < \infty$ implies that $0 = b^T A\bar{x} + b_0 F(A\bar{x}) \le \beta$, so that $\beta = 0$. Furthermore, we know that one of the inequalities above has to be strict for at least one element of $\Omega_1 \cup \Omega_2$, but this cannot be true on $\Omega_2$, so there exists $z \in \text{dom}(F)$ such that $b^T z < 0$. Since $b^T A\bar{x} = 0$ and $A\bar{x} \in \text{ridom}(F)$, we have $A\bar{x} - \epsilon(z - A\bar{x}) \in \text{dom}(F)$, so that $b^T(-\epsilon z) \le 0$, yielding a contradiction.

So, we need $b_0 \ne 0$, and the first pair of inequalities clearly requires $b_0 < 0$, so that we can take $b_0 = -1$. This shows that

$$b^T z - \beta \le F(z)$$

for all $z$ and

$$b^T Ay - \beta \ge g^T(y-x) + F(Ax)$$

for all $y$. Taking $y = x$, $z = Ax$, we find that $\beta = b^T Ax - F(Ax)$ yielding

$$F(z) - F(Ax) \ge b^T(z-x)$$

for all $z$ and $b^T A(y-x) \ge g^T(y-x)$ for all $y$. This last inequality implies that $g = A^T b$ and the first one that $b \in \partial F(Ax)$, therefore concluding the proof.

# Chapter 4

# Introduction: Bias, Variance and Density Estimation

In this chapter, we illustrate the bias variance dilemma in the context of density estimation, in which problems are similar to those encountered in classical parametric or non-parametric statistics [159, 60, 154].

For density estimation, one assumes that a random variable $X$ is given with unknown p.d.f. $f$ and we want to build an estimator, i.e., a mapping $(x, T) \mapsto \hat{f}(x; T)$ that provides an estimation of $f(x)$ based on a training set $T = (x_1, \ldots, x_N)$ containing $N$ i.i.d. realizations of $X$ (i.e., $T$ is a realization of $\mathbb{T} = (X_1, \ldots, X_N)$, $N$ independent copies of $X$). Alternatively, we will say that the mapping $T \mapsto \hat{f}(\cdot; T)$ is an estimator of the full density $f$. Note that, to further illustrate our notation, $\hat{f}(x; T)$ is a number while $\hat{f}(x; \mathbb{T})$ is a random variable.

## 4.1 Parameter estimation and sieves

Parameter estimation is the most common density estimation method, in which one restrict $\hat{f}$ to belong to a finite-dimensional parametric class, denoted $(f_\theta, \theta \in \Theta)$, with $\Theta \subset \mathbb{R}^p$. For example, $f_\theta$ can be a family of Gaussian distributions on $\mathbb{R}^d$. With our notation, a parametric model provides estimators taking the form

$$\hat{f}(x; T) = f_{\hat{\theta}(T)}(x)$$

and the problem becomes the estimation of the parameter $\hat{\theta}$.

There are several, well-known methods for parameter estimation, and, since this is not the focus of the book, we only consider the most common one, maximum

likelihood, which consists in computing $\hat{\theta}$ that maximizes the log-likelihood

$$C(\theta) = \frac{1}{N} \sum_{k=1}^{N} \log f_\theta(x_k). \tag{4.1}$$

The resulting $\hat{\theta}$ (when it exists) is called the maximum likelihood estimator of $\theta$, or m.l.e.

   If the true $f$ belongs to the parametric class, so that $f = f_{\theta_*}$ for some $\theta_* \in \Theta$, standard results in mathematical statistics [29, 118] provide sufficient conditions for $\hat{\theta}$ to converge to $\theta_*$ when $N$ tends to infinity. However, the fact that the true p.d.f. belongs to the finite dimensional class $(f_\theta)$ is an optimistic assumption that is generally false. In this regard, the standard theorems in parametric statistics may be regarded as analyzing a "best case scenario," or as performing a "sanity check," in which one asks whether, in the ideal situation in which $f$ actually belongs to the parametric class, the designed estimator has a proper behavior. In non-parametric statistics, a parametric model can still be a plausible approach in order to approximate the true $f$, but the relevant question should then be whether $\hat{f}$ provides (asymptotically), the best approximation to $f$ among all $f_\theta, \theta \in \Theta$. The maximum likelihood estimator can be analyzed from this viewpoint, if one measures the difference between two density functions by the Kullback-Liebler divergence (also called differential entropy):

$$KL(f\|f_\theta) = \int_{\mathbb{R}^d} \log \frac{f(x)}{f_\theta(x)} f(x)dx \tag{4.2}$$

which is positive unless $f = f_\theta$ (and may be equal to $+\infty$).

   This expression of the divergence is a simplification of its general measure-theoretic definition, that we now provide for completeness—and future use. Let $\mu$ and $\nu$ be two probability measures on a set $\widetilde{\Omega}$. One says that $\mu$ is absolutely continuous with respect to $\nu$, with notation $\mu \ll \nu$, if, for every (measurable) subset $A \subset \widetilde{\Omega}$, $\nu(A) = 0$ implies $\mu(A) = 0$. The Radon-Nikodym theorem [31] states that $\mu \ll \nu$ is and only if there exists a non-negative function $g$ defined on $\widetilde{\Omega}$ such that

$$\mu(A) = \int_A g(x)d\nu(x).$$

In terms of random variables, this says that, if $X : \Omega \to \widetilde{\Omega}$ and $Y : \Omega \to \widetilde{\Omega}$ are two random variables with respective distributions $\mu$ and $\nu$, and $\varphi : \widetilde{\Omega} \to \mathbb{R}$ is measurable, then $E(\varphi(X)) = E(g(Y)\varphi(Y))$. The function $g$ is called the Radon-Nikodym derivative of $\mu$ with respect to $\nu$ and is denoted $d\mu/d\nu$ (it is defined up to a modification on a set of $\nu$-probability zero). The general definition of the Kullback-Liebler divergence

between $\mu$ and $\nu$ is then:

$$KL(\mu\|\nu) = \begin{cases} \int_{\tilde{\Omega}} \left(\log \frac{d\mu}{d\nu}\right) \frac{d\mu}{d\nu} d\nu & \text{if } \mu \ll \nu \\ +\infty & \text{otherwise} \end{cases} \tag{4.3}$$

In the case when $\mu = f \, dx$ and $\nu = \tilde{f} \, dx$ are both probability measures on $\mathbb{R}^d$ with respective p.d.f.'s $f$ and $\tilde{f}$, $\mu \ll \nu$ means that $f/\tilde{f}$ is well defined everywhere except on a set of $\nu$-probability zero. It is then equal to $d\mu/d\nu$. If $\mu \ll \nu$, we can therefore write

$$KL(\mu\|\nu) = \int_{\mathbb{R}^d} \frac{f(x)}{\tilde{f}(x)} \log\left(\frac{f(x)}{\tilde{f}(x)}\right) \tilde{f}(x) dx = \int_{\mathbb{R}^d} \log\left(\frac{f(x)}{\tilde{f}(x)}\right) f(x) dx$$

and we will make the abuse of notation of writing $KL(f\|\tilde{f})$ for $KL(f \, dx\|\tilde{f} \, dx)$, which gives the expression provided in (4.2).

The general definition also gives a simple expression when $\widetilde{\Omega}$ is a finite set, with

$$KL(\mu\|\nu) = \sum_{x \in \widetilde{\Omega}} \log \frac{\mu(x)}{\nu(x)} \mu(x),$$

that we will use later in these notes (if there exists $x$ such that $\mu(x) > 0$ and $\nu(x) = 0$, then $KL(\mu\|\nu) = \infty$). The most important property for us is that the Kullback-Liebler divergence can be used as a measure of discrepancy between two probability distribution, based on the following proposition.

**Proposition 4.1** *Let $\mu$ and $\nu$ be two probability measures on $\widetilde{\Omega}$. Then $KL(\mu\|\nu) \geq 0$ and vanishes if and only if $\mu = \nu$.*

Proof Assume that $\mu \ll \nu$ since the statement is obvious otherwise and let $g = d\mu/d\nu$. We have $\int_{\widetilde{\Omega}} g d\nu = 1$ (since, by definition, it is equal to $\mu(\widetilde{\Omega})$) so that

$$KL(\mu\|\nu) = \int_{\widetilde{\Omega}} (g \log g + 1 - g) d\nu.$$

We have $t \log t + 1 - t \geq 0$ with equality if and only $t = 1$ (the proof being left to the reader) so that $KL(\mu\|\nu) = 0$ if and only if $g = 1$ with $\nu$-probability one, i.e., if and only if $\mu = \nu$. ∎

Minimizing $KL(f\|f_\theta)$ with respect to $\theta$ is equivalent to maximizing

$$E_f(\log f_\theta) = \int_{\mathbb{R}^d} \log f_\theta(x) f(x) dx,$$

and an empirical evaluation of this expectation is $\frac{1}{N} \sum_{k=1}^{N} \log f_\theta(x_k)$, which provides the maximum likelihood method. Seen in this context, consistency of the maximum likelihood estimator states that this estimator almost surely converges to a best approximator of the true $f$ in the class $(f_\theta, \theta \in \Theta)$. More precisely, if one assumes that the function $\theta \mapsto \log f_\theta(x)$ is continuous[1] in $\theta$ for almost all $x$ and that, for all $\theta \in \Theta$, there exists a small enough $\delta > 0$ such that

$$\int_{\mathbb{R}^d} \left( \sup_{|\theta' - \theta| < \delta} \log f_{\theta'}(x) \right) f(x) \, dx < \infty$$

then, letting $\Theta_*$ denote the set of maximizers of $E_f(\log f_\theta)$, and assuming that it is not empty, the maximum likelihood estimator $\hat{\theta}_N$ is such that, for all $\epsilon > 0$ and all compact subsets $K \subset \Theta$,

$$\lim_{N \to \infty} \mathbb{P}\left( d(\hat{\theta}_N, \Theta_*) > \epsilon \text{ and } \hat{\theta}_N \in K \right) \to 0$$

where $d(\hat{\theta}_N, \Theta_*)$ is the Euclidean distance between $\hat{\theta}_N$ and the set $\Theta_*$. The interested reader can refer to Van der Vaart [194], Theorem 5.14, for a proof of this statement. Note that this assertion does not exclude the situation in which $\hat{\theta}_N$ goes to infinity (i.e., steps out of ever compact subset $K$ in $\Theta$), and the boundedness of the m.l.e. is either asserted from additional properties of the likelihood, or by simply restricting $\Theta$ to be a compact set.

If $\Theta_* = \{\theta_*\}$ and the m.l.e. almost surely converges to $\theta_*$, the speed of convergence can also be quantified by a central limit theorem (see Van der Vaart [194], Theorem 5.23) ensuring that, in standard cases $\sqrt{N}(\hat{\theta}_N - \theta_*)$ converges to a normal distribution.

Even though these results relate our present subject to classical parametric statistics, they are not sufficient for our purpose, because, when $f \neq f_{\theta_*}$, the convergence of the m.l.e. to the best approximator in $\Theta$ still leaves a gap in the estimation of $f$. This gap is often called the *bias* of the class $(f_\theta, \theta \in \Theta)$. One can reduce it by considering larger classes (e.g., with more dimensions), but the larger the class, the less accurate the estimation of the best approximator becomes for a fixed sample size (the estimator has a larger *variance*). This issue is known as the "bias vs. variance dilemma," and to address it, it is necessary to adjust the class $\Theta$ to the sample size in order to optimally balance the two types of error (and all non-parametric estimation methods have at least one mechanism that allows for this). When the "tuning parameter" is the dimension of $\Theta$, the overall approach is often referred to as the *method of sieves* [83, 80], in which the dimension of $\Theta$ is increased as a function of $N$ in a suitable way.

---

[1] Upper-semi continuous is sufficient.

Gaussian mixture models provide one of the most popular choices with the method of sieves. Modeling in this setting typically follows some variation of the following construction. Fix a sequence $(m_N, N \geq 1)$ and let

$$\Theta_N = \left\{ f : f(x) = \sum_{j=1}^{m_N} \alpha_j \frac{e^{-|x-\mu_j|^2/2\sigma^2}}{(2\pi\sigma^2)^{d/2}}, \right.$$

$$\left. \mu_1, \dots, \mu_{m_N} \in \mathbb{R}^d, \alpha_1 + \dots + \alpha_{m_N} = 1, \alpha_1, \dots, \alpha_{m_N} \in [0, +\infty), \sigma > 0 \right\}. \quad (4.4)$$

There are therefore $(d+1)m_N$ free parameters in $\Theta_N$. The integer $m_N$ allows one to adjust the dimension of $\Theta_N$ and therefore controls the bias-variance trade-off. If $m_N$ tends to infinity "slowly enough," the m.l.e. will converges (almost surely) to the true p.d.f. $f$ [80]. However, determining optimal sequences $N \to m_N$ remains a challenging and largely unsolved problem.

In practice the computation of the m.l.e. in this context uses an algorithm called *EM*, for expectation-maximization. This algorithm will be described later in chapter 16.

## 4.2   Kernel density estimation

Kernel density estimators [150, 177, 178] provide alternatives to the method of sieves. They also lend themselves to some analytical developments that provide elementary illustrations of the bias-variance dilemma.

Define a kernel function as a function $K : \mathbb{R}^d \to [0, +\infty)$ such that

$$\int_{\mathbb{R}^d} K(x)dx = 1, \quad \int_{\mathbb{R}^d} |x|K(x)\,dx < \infty, \quad \int_{\mathbb{R}^d} xK(x)\,dx = 0. \quad (4.5)$$

Note that the third equation is satisfied, in particular, when $K$ is an even function, i.e., $K(-x) = K(x)$.

Given $K$ and a scalar $\sigma > 0$, the rescaled kernel is defined by

$$K_\sigma(x) = \frac{1}{\sigma^d} K\left(\frac{x}{\sigma}\right).$$

Using the change of variable $y = x/\sigma$ (so that $dy = dx/\sigma^d$) one sees that $K_\sigma$ satisfies (4.5) as soon as $K$ does.

Based on a training set $T = (x_1, \dots, x_N)$, the kernel density estimator defines the family of densities

$$\hat{f}_\sigma(x; T) = \frac{1}{N} \sum_{k=1}^{N} K_\sigma(x - x_k)$$

One has

$$\int_{\mathbb{R}^d} K_\sigma(x - x_k)\,dx = 1$$

so that it is clear that $\hat{f}_\sigma$ is a p.d.f. In addition,

$$\int_{\mathbb{R}^d} x K_\sigma(x - x_k)\,dx = \int_{\mathbb{R}^d} (y + x_k) K_\sigma(y)\,dy = x_k$$

so that

$$\int_{\mathbb{R}^d} x \hat{f}_\sigma(x; T)\,dx = \bar{x}$$

where $\bar{x} = (x_1 + \cdots + x_N)/N$.

A typical choice for $K$ is a Gaussian kernel, $K(y) = e^{-|y|^2/2}/(2\pi)^{d/2}$. In this case, the estimated density is a sum of bumps centered at the data points $x_1, \ldots, x_N$. The width of the bumps is controlled by the parameter $\sigma$. A small $\sigma$ implies less rigidity in the model, which will therefore be more affected by changes in the data: the estimated density will have a larger variance. The converse is true for large $\sigma$, at the cost of being less able to adapt to variations in the true density: the model has a larger bias (see fig. 4.1 and fig. 4.2).

As we now show, in order to get a consistent estimator, one needs to let $\sigma = \sigma_N$ depend on the size of the training set. We have, taking expectations with respect to training data,

$$
\begin{aligned}
E(\hat{f}_\sigma(x; T)) &= \frac{1}{N\sigma^d} \sum_{k=1}^{N} E\Big(K((x - X_k)/\sigma)\Big) \\
&= \frac{1}{\sigma^d} \int_{\mathbb{R}^d} K((x - y)/\sigma) f(y)\,dy \\
&= \int_{\mathbb{R}^d} K(z) f(x - \sigma z)\,dz
\end{aligned}
$$

The bias of the estimator, i.e., the average difference between $\hat{f}_\sigma(x; T)$ and $f(x)$ is therefore given by

$$E(\hat{f}_\sigma(x; T)) - f(x) = \int_{\mathbb{R}^d} K(z)(f(x - \sigma z) - f(x))\,dz.$$

Interestingly, this bias does not depend on $N$, but only on $\sigma$, and it is clear that, under mild continuity assumptions on $f$, it will go to zero with $\sigma$.

The variance of $\hat{f}_\sigma(x; T)$ is given by

$$\mathrm{var}(\hat{f}_\sigma(x; T)) = \frac{1}{N\sigma^{2d}} \mathrm{var}(K((x - X)/\sigma))$$

Figure 4.1: Kernel density estimators using a Gaussian kernel and various values of $\sigma$ when the true distribution of the data is a standard Gaussian (Orange: true density; Blue: estimated density, Red dots: training data).

Figure 4.2:  Kernel density estimators using a Gaussian kernel and various values of $\sigma$ when the true distribution of the data is a Gamma distribution with parameter 2 (Orange:  true density; Blue: estimated density, Red dots: training data).

with

$$\frac{1}{N\sigma^{2d}}\mathrm{var}(K((x-X)/\sigma)) \;=\; \frac{1}{N\sigma^{2d}}\int_{\mathbb{R}}^{d} K((x-y)/\sigma)^2 f(y)dy$$

$$-\frac{1}{N\sigma^{2d}}\left(\int_{\mathbb{R}^d} K((x-y)/\sigma)f(y)dy\right)^2$$

$$=\; \frac{1}{N\sigma^{d}}\int_{\mathbb{R}^d} K(z)^2 f(x-\sigma z)dz - \frac{1}{N}\left(\int_{\mathbb{R}^d} K(z)f(x-\sigma z)dz\right)^2$$

The total mean-square error of the estimator is

$$\mathbb{E}((\hat{f}_\sigma(x)-f(x))^2) = \mathrm{var}(\hat{f}_\sigma(x)) + (\mathbb{E}(\hat{f}_\sigma(x))-f(x))^2.$$

Clearly, this error cannot go to zero unless we allow $\sigma = \sigma_N$ to depend on $N$. For the bias term to go to zero, we know that we need $\sigma_N \to 0$, in which case we can expect the second term in the variance to decrease like $1/N$, while, for the first term to go to zero, we need $N\sigma_N^d$ to go to infinity. This illustrates the bias-variance dilemma: $\sigma_N$ must go to zero in order to cancel the bias, but not too fast in order to also cancel the variance. There is, for each $N$, an optimal value of $\sigma$ that minimizes the error, and we now proceed to a more detailed analysis and make this statement a little more precise.

Let us make a Taylor expansion of both bias and variance, assuming that $f$ has at least three bounded derivatives and that $\int_{\mathbb{R}^d}|x|^3 K(x)\,dx < \infty$. We can write

$$f(x-\sigma z) = f(x) - \sigma z^T\nabla f(x) + \frac{\sigma^2}{2}z^T\nabla^2 f(x)z + O(\sigma^3|z|^3),$$

where $\nabla^2 f(x)$ denotes the matrix of second derivatives of $f$ at $x$. Since $\int zK(z)dz = 0$, this gives

$$\mathbb{E}(\hat{f}_\sigma(x;\mathbb{T})) - f(x) = \frac{\sigma^2}{2}M_f(x) + o(\sigma^2)$$

with $M_f = \int K(z)z^T\nabla^2 f(x)z\,dz$. Similarly, letting $S = \int K^2(z)\,dz$,

$$\mathrm{var}(\hat{f}_\sigma(x)) = \frac{1}{N\sigma^d}\left(Sf(x) + o(\sigma^d + \sigma^2)\right).$$

Assuming that $f(x) > 0$, we can obtain an asymptotically optimal value for $\sigma$ by minimizing the leading terms of the mean square error, namely

$$\frac{\sigma^4}{4}M_f^2 + \frac{S}{N\sigma^d}f(x)$$

which yields $\sigma_N = O(N^{-1/(d+4)})$ and

$$\mathbb{E}((\hat{f}_{\sigma_N}(x;\mathbb{T}) - f(x))^2) = O(N^{-4/(d+4)}).$$

If $f$ has $r + 1$ derivatives, and $K$ has $r - 1$ vanishing moments (this excludes the Gaussian kernel) one can reduce this error to $N^{-\frac{2r}{2r+d}}$. These rates can be shown to be "optimal," in the "min-max" sense, which roughly expresses the fact that, for any other estimator, there exists a function $f$ for which the convergence speed is at least as "bad" as the one obtained for kernel density estimation.

This result says that, in order to obtain a given accuracy $\epsilon$ in the worst case scenario, $N$ should be chosen of order $(1/\epsilon)^{1+(d/2r)}$ which grows exponentially fast with the dimension. This is the *curse of dimensionality* which essentially states that the issue of density estimation may be intractable in large dimensions. The same statement is true also for most other types of machine learning problems. Since machine learning essentially deals with high-dimensional data, this issue can be problematic.

Obviously, because the min-max theory is a worst-case analysis, not all situations will be intractable for a given estimator, and some cases that are challenging for one of them may be quite simple for others: even though all estimators are "cursed," the way each of them is cursed differs. Moreover, while many estimators are optimal in the min-max sense, this theory does not give any information on "how often" an estimator performs better than its worst case, or how it will perform on a given class of problems. (For kernel density estimation, however, what we found was almost universal with respect to the unknown density $f$, which indicates that this estimator is not a good choice in large dimensions.)

Another important point with this curse of dimensionality is that data may very often appear to be high dimensional while it has a simple, low-dimensional structure, maybe because many dimensions are irrelevant to the problem (they contain, for example, just random noise), or because the data is supported by a non-linear low-dimensional space, such as a curve or a surface. This information is, of course, not available to the analysis, but can sometimes be inferred using some of the dimension reduction methods that will be discussed later in chapter 20. Sometimes, and this is also important, information on the data structure can be provided by domain knowledge, that is, by elements, provided by experts, that specify how the data has been generated (such as underlying equations) and reasonable hypotheses that are made in the field. This source of information should never be ignored in practice.

# Chapter 5

# Prediction: Basic Concepts

## 5.1 General Setting

The goal of prediction is to learn, based on training data, an input-output relationship between two random variables $X$ and $Y$, in the sense of finding, for a specified criterion, the best function of the input $X$ that predicts the output $Y$. (In statistics, $Y$ is often called the dependent variable, and $X$ the independent variable.) We will, as always, assume that all the variables mentioned in this chapter are defined on a fixed probability space $(\Omega, \mathbb{P})$. We assume that $X : \Omega \to \mathcal{R}_X$, where $\mathcal{R}_X$ is the input space, and $Y : \Omega \to \mathcal{R}_Y$, where $\mathcal{R}_Y$ is the output space. The input-output relationship is therefore captured by an unknown function $f : \mathcal{R}_X \to \mathcal{R}_Y$, the predictor.

The following two subclasses of prediction problems are important enough to have learned their own names and specific literature.

- Quantitative output: $\mathcal{R}_Y = \mathbb{R}^q$ (often with $q = 1$). One then speaks of a *regression problem*.

- Categorical output: $\mathcal{R}_Y = \{g_1, \ldots, g_q\}$ is a finite set. One then speaks of a *classification problem*.

In most cases, the input space is Euclidean, i.e., $\mathcal{R}_X = \mathbb{R}^d$. Note also that, in classification, instead of a function $f : \mathcal{R} \to \mathcal{R}_Y$, one sometimes estimates a function $f : \mathcal{R}_X \to \Pi(\mathcal{R}_Y)$, where $\Pi(\mathcal{R}_Y)$ is the space of probability distributions on $\mathcal{R}_Y$. We will return to this in remark 5.4.

The quality of a prediction is assessed through the definition of a *risk function*. Such a function, denoted $r$, is defined on $\mathcal{R}_Y \times \mathcal{R}_Y$, takes values in $[0, +\infty)$ and should be understood as

$$r(\text{True output}, \text{Predicted output}), \tag{5.1}$$

so that $r(y, y')$ assigns a cost to the situation in which a true $y$ is predicted by $y'$. Note that this definition is asymmetric, and there is no requirement that $r(y, y') = r(y', y)$. It is important to remember our convention that the first variable is the true observation and the second one is a place-holder for a prediction. Risk functions are also called *loss functions*, or simply cost functions and we will use these terms as synonyms.

The goal in prediction is to minimize the *expected risk*, also called the *generalization error*:

$$R(f) = \mathbb{E}(r(Y, f(X))).$$

We will prove that an optimal $f$ can be easily described based on the joint distribution of $X$ and $Y$ (which is, unfortunately, never available). We will need for this to use conditional expectations and conditional probabilities and proceed first to a reminder of their definitions and properties.

## 5.2   Conditional expectation

If $\xi : \Omega \to \mathcal{R}_\xi$ and $\eta : \Omega \to \mathcal{R}_\eta \subset \mathbb{R}^d$ are discrete random variables, then

$$\mathbb{P}(\eta = \eta \mid \xi = \xi) = \mathbb{P}(\eta = \eta, \xi = \xi)/\mathbb{P}(\xi = \xi)$$

if $\mathbb{P}(\xi = \xi) > 0$ and is undefined otherwise. Then, if $\eta$ is real-valued and discrete, one defines the conditional expectation of $\eta$ given $\xi$, denoted $\mathbb{E}(\eta \mid \xi)$, by

$$\mathbb{E}(\eta \mid \xi)(\omega) = \sum_{\eta \in \mathcal{R}_\eta} \eta \mathbb{P}(\eta = \eta \mid \xi = \xi(\omega))$$

for all $\omega$ such that $\mathbb{P}(\xi = \xi(\omega)) > 0$. Note that $\mathbb{E}(\eta \mid \xi)$ is a random variable, defined over $\Omega$. It however only depends on the values of $\xi$, in the sense that $\mathbb{E}(\eta \mid \xi)(\omega) = \mathbb{E}(\eta \mid \xi)(\omega')$ if $\xi(\omega) = \xi(\omega')$. We will use the notation

$$\mathbb{E}(\eta \mid \xi = \xi) = \sum_{\eta \in \mathcal{R}_\eta} \eta \mathbb{P}(\eta = \eta \mid \xi = \xi),$$

which is now a function defined on $\mathcal{R}_\xi$. One has $\mathbb{E}(\eta \mid \xi)(\omega) = \mathbb{E}(\eta \mid \xi = \xi(\omega))$.

One can characterize $\mathbb{E}(\eta \mid \xi)$ by the properties

$$\begin{cases} \mathbb{E}(\eta \mid \xi) \text{ is a function of } \xi \\ \forall f : \mathcal{R}_\eta \to \mathbb{R}, \mathbb{E}(\mathbb{E}(\eta \mid \xi)f(\xi)) = \mathbb{E}(\eta f(\xi)). \end{cases} \tag{5.2}$$

The proof that our definition of $\mathbb{E}(\eta \mid \xi)$ for discrete random variables is the only one satisfying these properties is left to the reader. The interest of reformulating the

definition of the conditional expectation via (5.2) is that this provides a definition that works for general random variables (with the additional assumption that $f$ is measurable), not only for discrete ones. We assume below that $(\mathcal{R}_\xi, \mathcal{S}_\xi)$ and $(\mathcal{R}_\eta, \mathcal{S}_\eta)$ are measurable spaces.

**Definition 5.1** *Assume that $\mathcal{R}_\eta = \mathbb{R}^d$. Let $\xi : \Omega \to \mathcal{R}_\xi$ and $\eta : \Omega \to \mathcal{R}_\eta$ be two random variables with $\mathbb{E}(|\eta|) < \infty$. The conditional expectation of $\eta$ given $\xi$ is a random variable $\zeta : \Omega \to \mathcal{R}_\eta$ such that*

*(i)  There exists a function $h : \mathcal{R}_\xi \to \mathbb{R}$ such that $\zeta = h \circ \xi$ almost surely.*

*(ii)  For any measurable function $g : \mathcal{R}_\xi \to [0, +\infty)$, one has*

$$\mathbb{E}(\eta g \circ \xi) = \mathbb{E}(\zeta g \circ \xi).$$

*The variable $\zeta$ is then denoted $\mathbb{E}(\eta | \xi)$ and the function $h$ in (i) is denoted $\mathbb{E}(\eta | \xi = \cdot)$.*

Importantly, functions $\zeta$ satisfying conditions (i) and (ii) always exists and are almost surely unique, in the sense that if another function $\zeta'$ satisfies these conditions, then $\zeta = \zeta'$ with probability one. One obtains an equivalent definition if one restricts functions $g$ in (ii) to indicators of measurable sets, yielding the condition that, if $A \subset \mathcal{R}_\xi$ is measurable,

$$\mathbb{E}(\eta \mathbf{1}_{\xi \in B}) = \mathbb{E}(\zeta \mathbf{1}_{\xi \in B}).$$

Taking $g(\xi) = 1$ for all $\xi \in \mathcal{R}_\xi$ in condition (ii), one gets the well-known identity

$$\mathbb{E}(\mathbb{E}(\eta | \xi)) = \mathbb{E}(\eta).$$

Moreover, for any function $g$ defined on $\mathcal{R}_\xi$ we have $\mathbb{E}(\eta g \circ \xi | \xi) = (g \circ \xi)\mathbb{E}(\eta | \xi)$, which can be checked by proving that the right-hand side satisfies the conditions (i) and (ii).

Conditional expectations share many of the properties of simple expectations. For example, if $\eta \le \eta'$, both taking scalar values, then $\mathbb{E}(\eta | \xi) \le \mathbb{E}(\eta' | \xi)$ almost surely. Jensen's inequality also holds: if $\gamma : \mathbb{R}^d \to \mathbb{R}$ is convex and $\gamma \circ \eta$ is integrable, then

$$\gamma \circ \mathbb{E}(\eta | \xi) \le \mathbb{E}(\gamma \circ \eta | \xi).$$

We will discuss convex functions in chapter 3, but two important examples for this section are $\gamma(\eta) = |\eta|$ and $\gamma(\eta) = |\eta|^2$. The first one implies that $|\mathbb{E}(\eta | \xi)| \le \mathbb{E}(|\eta| \mid bfxi)$ and, taking expectations on both sides: $\mathbb{E}(|\mathbb{E}(\eta | \xi)|) \le \mathbb{E}(|\eta|)$, the upper bound being finite by assumption. For the square norm, we find that, if $\eta$ is square integrable, then so is $\mathbb{E}(\eta | \xi)$ and

$$\mathbb{E}(|\mathbb{E}(\eta | \xi)|^2) \le \mathbb{E}(|\eta|^2).$$

If $\eta$ is square integrable, then this inequality shows that $\mathbb{E}(\eta \mid \xi)$ minimizes $\mathbb{E}[|\eta - \zeta|^2]$ among all square integrable functions $\zeta : \Omega \to \mathcal{R}_\eta$ that satisfy (i). In other terms, the conditional expectation is the optimal least-square approximation of $\eta$ by a function of $\xi$. To see this, just write

$$
\begin{aligned}
\mathbb{E}[|\eta - \zeta|^2 \mid \xi] &= \mathbb{E}[|\eta|^2 \mid \xi] - 2\mathbb{E}[\eta^T \zeta \mid \xi] + |\zeta|^2 \\
&= \mathbb{E}[|\eta|^2 \mid \xi] - 2\mathbb{E}(\eta \mid \xi)^T \zeta + |\zeta|^2 \\
&= \mathbb{E}[|\eta|^2 \mid \xi] - |\mathbb{E}(\eta \mid \xi)|^2 + |\mathbb{E}(\eta \mid \xi) - \zeta|^2 \\
&= \mathbb{E}[|\eta - \mathbb{E}(\eta \mid \xi)|^2 \mid \xi] + |\mathbb{E}(\eta \mid \xi) - \zeta|^2 \\
&\geq \mathbb{E}[|\eta - \mathbb{E}(\eta \mid \xi)|^2 \mid \xi]
\end{aligned}
$$

and taking expectations on both sides yields the desired result.

If $A$ is a measurable subset of $\mathcal{R}_\eta$, the conditional expectation $\mathbb{E}(\mathbf{1}_A \mid \xi)$ (resp. $\mathbb{E}(\mathbf{1}_A \mid \xi = \xi)$) is denoted $\mathbb{P}(\eta \in A \mid \xi)$ (resp. $\mathbb{P}(\eta \in A \mid \xi = \xi)$), or $P_\eta(A \mid \xi)$ (resp. $P_\eta(A \mid \xi = \xi)$). While these functions are defined separately up to modifications on sets of probability zero. Under general assumptions on the set $\mathcal{R}_\eta$ and its $\sigma$-algebra (always satisfied in our discussions), these conditional probabilities can be defined together so that, for all $\omega \in \Omega$ $A \mapsto P_\eta(A \mid \xi)(\omega)$ is a probability distribution on $\mathcal{R}_\eta$ such that

$$
\mathbb{E}(\eta \mid \xi) = \int_{\mathcal{R}_\eta} \eta \, \mathbb{P}(d\eta \mid \xi).
$$

Assume that the the sets $\mathcal{R}_\xi$ and $\mathcal{R}_\eta$ are equipped with measures, say $\mu_\xi$ and $\mu_\eta$ such that the joint distribution of $(\xi, \eta)$ is absolutely continuous with respect to $\mu\xi \otimes \mu_\eta$, so that there exists a function $\varphi : \mathcal{R}_\xi \times \mathcal{R}_\eta \to \mathbb{R}$ (the p.d.f. of $(\xi, \eta)$ with respect to $\mu\xi \otimes \mu_\eta$) such that

$$
\mathbb{P}(\xi \in A, \eta \in B) = \int_{A \times B} \varphi(\xi, \eta) \mu_\xi \otimes \mu_\eta(dx, d\eta).
$$

Then $P_\eta(\cdot \mid \xi)$ is absolutely continuous with respect to $\mu_\eta$, with density given by the conditional p.d.f. of $\eta$ given $\xi$, namely,

$$
\varphi(\cdot \mid \xi) : (\eta, \omega) \mapsto \frac{\varphi(\eta, \xi(\omega))}{\int_{\mathcal{R}_\eta} \varphi(\eta', \xi(\omega)) \mu_\eta(d\eta')} = \varphi(\eta \mid \xi = \xi(\omega)). \tag{5.3}
$$

Note that

$$
\mathbb{P}\left\{ \omega : \int_{\mathcal{R}_\eta} \varphi(\eta', \xi(\omega)) \mu_\eta(d\eta') = 0 \right\} = 0
$$

so that the conditional density can be defined arbitrarily when the numerator vanishes[1].

The most common example is when $\mathcal{R}_\xi$ and $\mathcal{R}_\eta$ are Euclidean spaces and $\mu_\xi$, $\mu_\eta$ are Lebesgue's measures, in which case (5.3) is the usual definition of conditional p.d.f.'s. Note also that, for discrete random variables, (5.3) coincides with the definition of conditional probabilities $\mathbb{P}(\eta = \eta \mid \xi = \xi(\omega))$ when $\mu_x$ and $\mu_\eta$ are counting measures. As a last example, if $\mathcal{R}_\xi = \mathbb{R}^d$, $\mu_\xi$ is Lebesgue's measure and $\eta$ is discrete, then

$$\varphi(\eta \mid \xi = \xi(\omega)) = \frac{\varphi(\eta, \xi(\omega))}{\sum_{\eta' \in \mathcal{R}_\eta} \varphi(\eta', \xi(\omega))}.$$

## 5.3   Bayes predictor

Recall that $r : (y, y') \mapsto r(y, y')$ denotes the risk function and that we want to minimize $R(f) = \mathbb{E}(r(Y, f(X)))$ over all possible predictors $f$.

**Definition 5.2** *A Bayes predictor is a measurable function $f : \mathcal{R}_X \to \mathcal{R}_Y$ such that, for all $x \in \mathcal{R}_X$,*

$$\mathbb{E}\Big(r(Y, f(x)) \mid X = x\Big) = \min\Big\{\mathbb{E}\Big(r(Y, y') \mid X = x\Big) : y' \in \mathcal{R}_Y\Big\}$$

There can be multiple Bayes predictors if the minimum in the proposition is not uniquely attained. Note that, if $f^*$ is a Bayes predictor and $\hat{f}$ any other predictor, we have, by definition

$$\mathbb{E}\Big(r(Y, f^*(X)) \mid X\Big) \le \mathbb{E}\Big(r(Y, \hat{f}(X)) \mid X\Big).$$

Passing to expectations, this implies $R(f^*) \le R(\hat{f})$. We therefore have the following result:

**Theorem 5.3** *Any Bayes predictor $f^*$ is optimal, in the sense that it minimizes the generalization error $R$.*

*Example 1. Regression with mean-square error.* When $\mathcal{R}_X = \mathbb{R}^d$ and $\mathcal{R}_Y = \mathbb{R}^q$, the most common risk function is the squared norm $r(y, y') = |y - y'|^2$. The resulting generalization error is called the MSE (mean square error) and given by $R(f) = \mathbb{E}(|Y - f(X)|^2)$. The Bayes predictor is such that $f^*(x)$ minimizes

$$t \mapsto \mathbb{E}(|Y - t|^2 \mid X = x).$$

---

[1]Letting $\varphi_\xi(\xi) = \int_{\mathcal{R}_\eta} \varphi(\eta', \xi) \mu_\eta(d\eta')$, which is the marginal p.d.f. of $\xi$ with respect to $\mu_\xi$, we have

$$\mathbb{P}(\varphi_\xi(\xi) = 0) = \int_{\mathcal{R}_\xi} \mathbf{1}_{\varphi_\xi(\xi) = 0} \varphi_\xi(\xi) \mu_\xi(d\xi) = 0.$$

Let $f^*(x) = E(Y \mid X = x)$ and write

$$
\begin{aligned}
\mathbb{E}(|Y - t|^2 \mid X = x) =& \mathbb{E}(|Y - f^*(x)|^2 \mid X = x) + 2\mathbb{E}((Y - f^*(x))^T(f^*(x) - t) \mid X = x) \\
& + |f^*(x) - t|^2 \\
=& \mathbb{E}(|Y - f^*(x)|^2 \mid X = x) + 2\mathbb{E}((Y - f^*(x))^T \mid X = x)(f^*(x) - t) \\
& + |f^*(x) - t|^2 \\
=& \mathbb{E}(|Y - f^*(x)|^2 \mid X = x) + |f^*(x) - t|^2.
\end{aligned}
$$

This proves that $\mathbb{E}(Y \mid X = x)$ is the unique Bayes classifier (up to a modification on a set of probability 0).

*Example 2. Classification with zero-one loss.* Let $\mathcal{R}_X = \mathbb{R}^d$ and $\mathcal{R}_Y$ be a finite set. The zero-one loss function is defined by $r(y, y') = 1$ if $y \neq y'$ and $0$ otherwise. From this, it results that the generalization error is the probability of misclassification $R(f) = P(Y \neq f(X))$ (also called the misclassification error).

The Bayes predictor is such that $f^*(x)$ minimizes

$$
g \mapsto \mathbb{P}(Y \neq g \mid X = x) = 1 - \mathbb{P}(Y = g \mid X = x).
$$

It is therefore given by the so-called *posterior mode*:

$$
f^*(x) = \operatorname{argmax}_g \mathbb{P}(Y = g \mid X = x).
$$

**Remark 5.4** As mentioned at the beginning of the chapter, one sometimes replaces a pointwise prediction of the output by a probabilistic one, so that $f(x)$ is a probability distribution on $\mathcal{R}_Y$. If $A$ is a (measurable) subset of $\mathcal{R}_Y$, we will write $f(x, A)$ rather than $f(x)(A)$.

In such a case, the loss function, $r$, is defined on $\mathcal{R}_Y \times \Pi(\mathcal{R}_Y)$, and the expected risk is still defined by $\mathbb{E}(r(Y, f(X)))$.

It is quite natural to require that $\pi \mapsto r(y, \pi)$ is minimized. For classification problems, where $\mathcal{R}_Y$ is finite, one can choose

$$
r(y, \pi) = -\log \pi(y) \tag{5.4}
$$

The Bayes estimator is then a minimizer of $\pi \mapsto -\mathbb{E}(\log \pi(Y) \mid X = x)$. The solution is (unsurprisingly) $f(x, y) = \mathbb{P}(Y = y \mid X = x)$ since we always have

$$
-\mathbb{E}(\log \pi(Y) \mid X = x) = -\sum_{y \in \mathcal{R}_Y} \log \pi(y) f(x, y) \geq -\sum_{y \in \mathcal{R}_Y} \log f(x, y) f(x, y).
$$

The difference between these terms is indeed

$$
\sum_{y \in \mathcal{R}_Y} \log \frac{f(x, y)}{\pi(y)} f(x, y) = KL(f(x, \cdot), \pi) \geq 0.
$$

For regression problems, with $\mathcal{R}_Y = \mathbb{R}^q$, one can choose

$$r(y, \pi) = \int_{\mathbb{R}^q} |z - y|^2 \pi(dz)$$

which is indeed minimum when $\pi$ is concentrated on $y$. Here, the Bayes estimator minimizes (with respect to $\pi$)

$$\int_{\mathbb{R}^q} \int_{\mathbb{R}^q} |z - y|^2 \pi(dz) P_Y(x, dy) = \int_{\mathbb{R}^q} \left( \int_{\mathbb{R}^q} |z - y|^2 P_Y(x, dy) \right) \pi(dz)$$

where $P_Y(x, \cdot)$ is the conditional distribution of $Y$ given $X = x$. For any $z$, one has

$$\int_{\mathbb{R}^q} |y - z|^2 f(x, dy) \geq \int_{\mathbb{R}^q} |y - \mathbb{E}(Y \mid X = x)|^2 f(x, dy) \qquad \blacklozenge$$

which shows that the Bayes estimator is, in this case, the Dirac measure concentrated at $\mathbb{E}(Y \mid X = x)$.

## 5.4 Examples: model-based approach

Bayes predictors are never available in practice, because the true distribution of $(X, Y)$, or that of $Y$ given $X$, are unknown. These distributions can only be inferred from observations, i.e., from a training set: $T = (x_1, y_1, \ldots, x_N, y_N)$.

This is the approach followed by *model-based*, or *generative* methods, namely using training data to approximate the joint distribution of $X$ and $Y$ with a statistical model estimated from data before using the Bayes estimator derived from this model for prediction. We now illustrate this approach with a few examples.

### 5.4.1 Gaussian models and naive Bayes

Consider a regression problem with $\mathcal{R}_Y = \mathbb{R}$, and model the joint distribution of $(X, Y)$ as a $(d + 1)$-dimensional Gaussian distribution with mean $\mu$ and covariance matrix $\Sigma$, which must be estimated from data. Write $\mu = \begin{pmatrix} m \\ \mu_0 \end{pmatrix}$, with $\mu_0 \in \mathbb{R}$, $m \in \mathbb{R}^d$ and $\Sigma$ in the form, for some symmetric matrix $S$ and $d$-dimensional vector $u$

$$\Sigma = \begin{pmatrix} S & u \\ u^T & \sigma_{00}^2 \end{pmatrix}.$$

Then, letting $\Delta = \sigma_{00}^2 - u^T S^{-1} u$,

$$\Sigma^{-1} = \frac{1}{\Delta} \begin{pmatrix} \Delta S^{-1} + S^{-1} u u^T S^{-1} & -S^{-1} u \\ -u^T S^{-1} & 1 \end{pmatrix}.$$

This shows that the joint p.d.f. of $(X, Y)$ is proportional to

$$\exp\left(-\frac{1}{2\Delta}\left((y - \mu_0)^2 - 2u^T S^{-1}(x - m)(y - \mu_0) + (\textit{terms not depending on } y)\right)\right).$$

In particular

$$E(Y|X) = \mu_0 + u^T S^{-1}(x - m),$$

which provides the least-square linear regression predictor. (In this expression, $u$ is the covariance between $X$ and $Y$ and $S$ is the covariance matrix of $X$.)

If one restricts the model to having a diagonal covariance matrix $S$, then

$$E(Y|X) = \mu_0 + \sum_{j=1}^{d} \frac{u^{(j)}}{s_{jj}}(x^{(j)} - m^{(j)}).$$

This predictor is often called the *naive Bayes* predictor for regression.

### 5.4.2 Kernel regression

Let $\mathcal{R}_X = \mathbb{R}^d$ and $\mathcal{R}_Y = \mathbb{R}$. Let $K_1 : \mathbb{R}^d \to \mathbb{R}$ and $K_2 : \mathbb{R} \to \mathbb{R}$ be two kernels, therefore satisfying

$$\int_{\mathbb{R}^d} K_1(x)dx = \int_{\mathbb{R}} K_2(x)dx = 1; \quad \int_{\mathbb{R}^d} xK_1(x)dx = \int_{\mathbb{R}} yK_2(y)dy = 0.$$

Let $K(x,y) = K_1(x)K_2(y)$ so that

$$\int_{\mathbb{R}^{d+1}} K(x,y)dydx = 1$$

$$\int_{\mathbb{R}^{d+1}} yK(y,x)dydx = 0$$

$$\int_{\mathbb{R}^{d+1}} xK(y,x)dydx = 0.$$

The kernel estimator of the joint p.d.f., $\varphi$, of $(X, Y)$ at scale $\sigma$ is, in this case:

$$\hat{\varphi}(x,y) = \frac{1}{N} \sum_{k=1}^{N} \frac{1}{\sigma^{d+1}} K_1\left(\frac{x - x_k}{\sigma}\right) K_2\left(\frac{y - y_k}{\sigma}\right).$$

Based on $\hat{\varphi}$, the conditional expectation of $Y$ given $X$ is

$$\hat{f}(x) = \frac{\frac{1}{N} \sum_{k=1}^{N} \frac{1}{\sigma^{d+1}} \int_{\mathbb{R}} yK_1\left(\frac{x-x_k}{\sigma}\right) K_2\left(\frac{y-y_k}{\sigma}\right) dy}{\frac{1}{N} \sum_{k=1}^{N} \frac{1}{\sigma^{d+1}} \int_{\mathbb{R}} K_1\left(\frac{x-x_k}{\sigma}\right) K_2\left(\frac{y-y_k}{\sigma}\right) dy}.$$

Using the fact that $\sigma^{-1} \int_{\mathbb{R}} y K_2\left(\frac{y-y_k}{\sigma}\right) dy = y_k$, we can simplify this expression to obtain

$$\hat{f}(x) = \frac{\sum_{k=1}^{N} y_k K_1\left(\frac{x-x_k}{\sigma}\right)}{\sum_{k=1}^{N} K_1\left(\frac{x-x_k}{\sigma}\right)}.$$

This the kernel-density regression estimator [139, 202].

### 5.4.3   A classification example

Let $\mathcal{R}_Y = \{0, 1\}$ and assume $\mathcal{R}_X = \mathbb{N} = \{0, 1, 2, \ldots\}$. Let $p = \mathbb{P}(Y = 1)$ and assume that conditionally to $Y = g$, $X$ follows a Poisson distribution with mean $\mu_g$. Assume that $\mu_0 < \mu_1$.

The posterior distribution of $Y$ given $X = x$ is [2]

$$\mathbb{P}(Y = g \mid X = x) \propto \begin{cases} (1-p)\mu_0^x e^{-\mu_0} & \text{if } g = 0 \\ p\mu_1^x e^{-\mu_1} & \text{if } g = 1 \end{cases}$$

A Bayes classifier is then provided by taking $f(x) = 1$ if

$$\log p + x \log \mu_1 - \mu_1 \geq \log(1-p) + x \log \mu_0 - \mu_0$$

that is:

$$x \log \frac{\mu_1}{\mu_0} \geq \log \frac{1-p}{p} + \mu_1 - \mu_0$$

Since we are assuming that $\mu_1 > \mu_0$, we find that $f(x) = 1$ if [3]

$$x \geq \left\lceil \frac{\log((1-p)/p) + \mu_1 - \mu_0}{\log(\mu_1/\mu_0)} \right\rceil$$

and 0 otherwise.

## 5.5   Empirical risk minimization

### 5.5.1   General principles

Model-based approaches for prediction are based on the estimation of the joint distribution of the input and output variables, which is arguably a harder problem than prediction [196]. Since the goal is to find $f$ minimizing the expected risk

---

[2] $\propto$ is the notation for "proportional to"
[3] $\lceil x \rceil$ is the smallest integer larger than $x$ (ceiling).

$R(f) = \mathbb{E}(r(Y, f(X)))$, one may prefer a direct approach and consider the minimization of an empirical estimate of this risk, based on training data $T = (x_1, y_1, \ldots, x_N, y_N)$, namely

$$\hat{R}(f) = \frac{1}{N} \sum_{k=1}^{N} r(y_k, f(x_k)).$$

This strategy is called *empirical risk minimization.*

Importantly, $\hat{R}$ must be minimized over a restricted class, $\mathcal{F}$, of predictors to avoid overfitting. For example, with $\mathcal{R}_Y = \mathbb{R}$ and $\mathcal{R} = \mathbb{R}^d$, one can take

$$\mathcal{F} = \left\{ f : f(x) = \beta_0 + \sum_{i=1}^{d} b^{(i)} x^{(i)} : \ \beta_0, b^{(1)} \ldots, b^{(d)} \in \mathbb{R} \right\}.$$

Minimizing the empirical mean-square error

$$\hat{R}(f) = \frac{1}{N} \sum_{k=1}^{N} (y_k - f(x_k))^2$$

over $f \in \mathcal{F}$ leads to the standard least-square regression estimator.

As another example, consider

$$\mathcal{F} = \left\{ f : f(x) = \sum_{j=1}^{p} w_j \psi \left( \beta_{j0} + \sum_{i=1}^{d} \beta_{ji} x^{(i)} \right), w_j, \beta_{ji} \in \mathbb{R} \right\}.$$

with a fixed function $\psi$. This corresponds to a two-layer perceptron model.

As a last example for now (we will see many others in the rest of this book), taking $d = 1$, the set

$$\mathcal{F} = \left\{ f : \int_{\mathbb{R}} f''(x)^2 dx < \mu \right\}$$

(with $\mu > 0$) provides an infinite dimensional space of predictors, which leads to spline regression.

### 5.5.2  Bias and variance

We give a further illustration of the bias-variance dilemma in the regression case, using the mean-square error and taking $q = 1$ to simplify. Denote the Bayes predictor by $f^*(x) = \mathbb{E}(Y \mid X = x)$.

Fix a function space $\mathcal{F}$, and let $\hat{f}^*$ be the optimal predictor in $\mathcal{F}$, in the sense that it minimizes $E(|Y - f(X)|^2)$ over $f \in \mathcal{F}$. Then, letting $\hat{f}_N \in \mathcal{F}$ denote an estimated predictor,

$$
\begin{aligned}
R(\hat{f}_N) &= \mathbb{E}(|Y - \hat{f}_N(X)|^2) \\
&= \mathbb{E}(|Y - \hat{f}^*(X)|^2) + \mathbb{E}(|\hat{f}_N(X) - \hat{f}^*(X)|^2) \\
&\quad + 2\mathbb{E}((Y - \hat{f}^*(X))(\hat{f}^*(X) - \hat{f}_N(X))
\end{aligned}
$$

Let us make the assumption that there exists $\epsilon > 0$ such that $f_\lambda = \hat{f}^* + \lambda(\hat{f}_N - \hat{f}^*)$ belongs to $\mathcal{F}$ for $\lambda \in [-\epsilon, \epsilon]$. This happens when $\mathcal{F}$ is a linear space, or more generally when $\mathcal{F}$ is convex and $\hat{f}^*$ is in its relative interior (see chapter 3). Let $\psi : \lambda \mapsto \mathbb{E}(|Y - f_\lambda(X)|^2)$, which is minimal at $\lambda = 0$. We have

$$
\begin{aligned}
\psi(\lambda) &= \mathbb{E}(|Y - \hat{f}^*(X) - \lambda(\hat{f}_N(X) - \hat{f}^*(X))|^2) \\
&= \mathbb{E}(|Y - \hat{f}^*(X)|^2) - 2\lambda\mathbb{E}((Y - \hat{f}^*(X))(\hat{f}_N(X) - \hat{f}^*(X))) + \lambda^2\mathbb{E}(|\hat{f}_N(X) - \hat{f}^*(X)|^2)
\end{aligned}
$$

and

$$
0 = \psi'(0) = 2\mathbb{E}((Y - \hat{f}^*(X))(\hat{f}^*(X) - \hat{f}_N(X)))
$$

We therefore get the identity

$$
R(\hat{f}_N) = \mathbb{E}(|Y - \hat{f}^*(X)|^2) + \mathbb{E}(|\hat{f}_N(X) - \hat{f}^*(X)|^2) = \text{``Bias''} + \text{``Variance''}.
$$

The bias can be further decomposed as

$$
\mathbb{E}(|Y - \hat{f}^*(X)|^2) = \mathbb{E}(|Y - f^*(X)|^2) + \mathbb{E}(|f^*(X) - \hat{f}^*(X)|^2)
$$

because $f^*$ is the conditional expectation. As a result, we obtain an expression the generalization error with three contributions, namely,

$$
R(\hat{f}_N) \le \mathbb{E}(|Y - f^*(X)|^2) + \mathbb{E}(|f^* - \hat{f}^*(X)|^2) + \mathbb{E}(|\hat{f}_N(X) - \hat{f}^*(X)|^2).
$$

The first term is the Bayes error. It is fixed by the joint distribution of $X$ and $Y$ and measures how well $Y$ can be approximated by a function of $X$. The second term compares $f^*$ to its best approximation in $\mathcal{F}$, and is therefore reduced by taking larger model spaces. The last term is the error caused by using the data to estimate $\hat{f}^*$. It increases with the size of $\mathcal{F}$. This is illustrated in Figure 5.1.

**Remark 5.5** If the assumption made on $\hat{f}^*$ is not valid, one can write

$$
R(\hat{f}_N) = \mathbb{E}(|Y - \hat{f}_N(X)|^2) \le 2\left(\mathbb{E}(|Y - \hat{f}^*(X)|^2) + \mathbb{E}(|\hat{f}_N(X) - \hat{f}^*(X)|^2)\right)
$$

and still obtain a control (as an inequality) of the generalization error by a bias-plus-variance sum. ◆

Figure 5.1:  Sources of errors in statistical Learning: When $P^*$ is the distribution of the data, the optimal predictor $f^*$ minimizes the expected loss function.  Based on data $Z_1, \ldots, Z_N$, the sample-based distribution is $\hat{P} = (\delta_{Z_1} + \cdots + \delta_{Z_N})/N$ and the empirical loss is minimized over a subset $\mathcal{S}$ of the space of all possible estimators.  The expected discrepancy between the resulting estimator and the one minimizing the true expected loss on the subspace is the "variance" of the method, and the expected discrepancy between this subspace-constrained estimator and and the optimal one is the "bias."

## 5.6 Evaluating the error

### 5.6.1 Generalization error

Given input and output variables $X : \Omega \to \mathcal{R}_X$ and $Y : \Omega \to \mathcal{R}_Y$ and a risk function $r : \mathcal{R}_Y \times \mathcal{R}_Y \to [0. + \infty)$, we have defined the generalization (or prediction) error as

$$R(f) = \mathbb{E}(r(Y, f(X))).$$

Recall that a training set $T = ((x_1, y_1), \ldots, (x_N, y_N))$ is a realization $T = \mathbb{T}(\omega)$ of the random variable $\mathbb{T} = ((X_1, Y_1), \ldots, (X_N, Y_N))$, an i.i.d. sample of the joint distribution of $(X, Y)$. A *learning algorithm* is a function $T \mapsto \hat{f}_T$ defined on the set of training sets, namely, $\bigcup_{N=1}^{\infty} (\mathcal{R} \times \mathcal{R}_Y)^N$ and taking values in $\mathcal{F}$.

For a given $T$ and a specific algorithm, one is primarily interested in evaluating $R(\hat{f}_T)$, the generalization error of the predictor estimated from observed data. To emphasize the fact that the training set is fixed in this expression, one often writes:

$$R(\hat{f}_T) = \mathbb{E}(r(Y, \hat{f}_{\mathbb{T}}(X))|\mathbb{T} = T)$$

If we also take the expectation with respect to $T$ (for fixed $N$), we obtain the averaged generalization risk as

$$\mathbb{E}(R(\hat{f}_{\mathbb{T}})) = \mathbb{E}(r(Y, \hat{f}_{\mathbb{T}}(X))),$$

which provides an evaluation of the average quality of the algorithm when evaluated on random training sets of size $N$. If $A : T \mapsto \hat{f}_T$ denotes the learning algorithm, we will denote $\overline{R}_N(A) = \mathbb{E}(R(\hat{f}_{\mathbb{T}}))$.

Since their computation requires the knowledge of the joint distribution of $X$ and $Y$, these errors are not available in practice. Given a training set $T$ and a predictor $f$, one can compute the empirical error

$$\hat{R}_T(f) = \frac{1}{N} \sum_{k=1}^{N} r(y_k, f(x_k)).$$

Under the usual moment conditions, the law of large numbers implies that $\hat{R}_{\mathbb{T}}(f) \to R(f)$ with probability one for any given predictor $f$. However, the law of large numbers cannot be applied to assess whether the *in-sample error*,

$$\mathcal{E}_T \overset{\Delta}{=} \hat{R}_T(\hat{f}_T) = \frac{1}{N} \sum_{k=1}^{N} r(y_k, \hat{f}_T(x_k)),$$

is a good approximation of the generalization error $R(\hat{f}_T)$. This is because each term in the sum depends on the full data set, so that $\mathcal{E}_{\mathbb{T}}$ is not a sum of independent terms. The in-sample error typically under-estimates the generalization error, sometimes with a large discrepancy.

When one has enough data, however, it is possible to set some of it aside to form a test set. Formally, a test set is a collection $T' = (x'_1, y'_1, \ldots, x'_{N'}, y'_{N'})$ considered as a realization of an i.i.d. sample of $(X, Y)$, $\mathbb{T}' = (X'_1, Y'_1, \ldots, X'_{N'}, Y'_{N'})$, independent of $\mathbb{T}$. The test set error is then given by

$$\mathcal{E}_{T,T'} = \hat{R}_{T'}(\hat{f}_T) = \frac{1}{N'} \sum_{k=1}^{N'} r(y'_k, \hat{f}_T(x'_k)).$$

The law of large numbers (applied conditionally to $\mathbb{T} = T$) implies that $\mathcal{E}_{T,\mathbb{T}'}$ converges to $R(\hat{f}_T)$ with probability one when $N' \to \infty$.

However, in many applications, data acquisition is difficult or expensive (e.g., in the medical field) and sparing a part of it in order to form a test set is not a reasonable option. In such situations, cross-validation is generally a preferred alternative.

### 5.6.2 Cross validation

**Cross-validation error**

The $n$-fold cross-validation method (see, e.g., Stone [184]) separates the training set into $n$ non-overlapping sets of equal sizes, and estimates $n$ predictors by leaving out one of these subsets as a temporary test set. A generalization error is estimated from each test set and averaged over the $n$ results.

Let us formalize this computation after introducing some notation. We represent training data in the form $T = (z_1, \ldots, z_N)$, a sample of a random variable $Z$. With this notation, we can include supervised problems, such as prediction (taking $Z = (X, Y)$) and unsupervised ones such as density estimation (taking $Z = X$). One tries to estimate a function $f$ within a given class (e.g., a predictor, or a density) and one has a measure of "loss", denoted $\ell(f, z) \geq 0$ measuring how badly $f$ performs on the data $z$. For prediction, one takes $\ell(f, z) = r(y, f(x))$ with $z = (x, y)$ and for density estimation, e.g., $\ell(f, z) = -\log f(z)$, the negative log-likelihood. One then lets $R(f) = E(\ell(f, Z))$. For an algorithm $A : T \mapsto \hat{f}_T$, the loss $\bar{R}(A)$ is the quantity of interest.

Given another set $T' = (z'_1, \dots, z'_{N'})$, the empirical loss is

$$\hat{R}_{T'}(f) = \frac{1}{N'} \sum_{k=1}^{N'} \ell(f, z'_k)$$

and, using $T$ as a training set and $T'$ as a test set, we let

$$\mathcal{E}_{T,T'} = \hat{R}_{T'}(\hat{f}_T).$$

To define an $n$-fold cross-validation estimator of the error, one assumes that the training set $T$ is partitioned into $n$ subsets of equal sizes (up to one element if $N$ is not a multiple of $n$), $T_1, \dots, T_n$, so that $T_i$ and $T_j$ are non-intersecting if $i \ne j$, and $T = \bigcup_{i=1}^{n} T_i$. For each $i$, let $T^{(i)} = T \setminus T_i$, which provides the training data with the elements of $T_i$ removed. Then, the $n$-fold cross-validation error is defined by

$$\mathcal{E}_{\text{CV}}(T) = \frac{1}{n} \sum_{i=1}^{n} \mathcal{E}_{T^{(i)}, T_i}.$$

Assuming, to simplify, that $N$ is a multiple of $n$, the expectation of the cross-validation error is $E(R(\hat{f}_{\mathbb{T}_{N'}}))$, where the average is made over training sets $\mathbb{T}_{N'}$ of size $N' = N - N/n$. Note that the cross-validation error is an estimate of the average error of the algorithm over random training sets, not necessarily that of the current estimator $\hat{f}_T$. It returns an evaluation of the algorithm $A : T \mapsto \hat{f}_T$. When needed, one can emphasize this and write $\bar{R}_{\text{CV},T}(A)$.

The limit case when $n = N$ is called leave-one-out (LOO) cross validation. In this case $\mathcal{E}_{\text{CV}}$ is an almost unbiased estimator of $E(R(\hat{f}_T))$, but, because it is an average of functions of the training set that are quite similar (and that will therefore be positively correlated), its variance (as a function of $T$) may be quite large. Conversely, smaller values of $n$ will have smaller variances, but larger biases. In practice, it is difficult to assess which choice of $n$ is optimal, although 5- or 10-fold cross-validation is quite popular. LOO cross-validation is also often used, especially when $N$ is small.

**Model selection using cross validation**

Because it evaluates the quality of an algorithm, cross-validation is often used to perform model selection. Indeed, many learning algorithms depends on a parameter, that we will denote $\lambda$. In kernel density estimation, for example, $\lambda = \sigma$ is the kernel width. For mixtures of Gaussian, $\lambda = m$ is the number of Gaussian terms in the mixtures. Formally, this means that one has, for every $\lambda$, an algorithm $A_\lambda : T \mapsto \hat{f}_{T,\lambda}$.

Fixing a training set $T$, one can compute, for each $\lambda$, the cross-validation error $e_T(\lambda) = \bar{R}_{\mathrm{CV},T}(A_\lambda)$. Model selection is then performed by finding

$$\lambda^*(T) = \underset{\lambda}{\mathrm{argmin}}\, e_T(\lambda).$$

Once this $\lambda^*$ is obtained, the final estimator is $\hat{f}_{T,\lambda^*(T)}$, obtained by rerunning the algorithm one more time on the full training set.

This defines a new training algorithm, $A^* : T \mapsto \hat{f}_{T,\lambda^*(T)}$. It is a common mistake to consider that the cross-validation error associated to this algorithm is still given by $e(\lambda^*(T))$. This is false, because the computation of $\lambda^*$ uses the full training set. To compute the cross-validation error of $A^*$, one needs to encapsulate this model selection procedure in an other cross-validation loop. So, one needs to compute, using the previous notation,

$$\mathcal{E}^*_{\mathrm{CV}}(T) = \frac{1}{n} \sum_{i=1}^{n} \hat{R}_{T_i}(\hat{f}_{T^{(i)},\lambda^*(T^{(i)})})$$

where each $\hat{f}_{T^{(i)},\lambda^*(T^{(i)})}$ is computed by running a cross-validated model selection procedure restricted to $T^{(i)}$. This is often called a double-loop cross-validation procedure (the number of folds in the inner and outer loops do not have to coincide). Note that each $\lambda^*(T^{(i)})$ that does not necessarily coincide with the optimal $\lambda^*(T)$ obtained with the full training set.

# Chapter 6

# Inner Products and Reproducing Kernels

## 6.1  Introduction

We will discuss later in this book various methods that specify the prediction is as a linear function of the input. These methods are often applied after taking transformations of the original variables, in the form $x \mapsto h(x)$ (i.e., the prediction algorithm is applied to $h(x)$ instead of $x$). We will refer to $h$ as a "feature function," which typically maps the initial data $x \in \mathcal{R}$ to a vector space, sometimes of infinite dimensions, that we will denote $H$ (the "feature space").

The present chapter provides a formal description of this framework, focusing, in particular, on situations in which $H$ has an inner product, as this inner product is often instrumental in the design of linear methods on $H$. Many machine learning methods can indeed be expressed either as functions of the coordinates of the input data in some space, or as functions of the inner products between the input samples. Such methods can bypass the difficulty of using high-dimensional features with the help of the theory of "reproducing kernels," [12, 201] which ensures that the inner product between special classes of feature functions $h(x)$ and $h(x')$ can be explicitly computed as a function of $x$ and $x'$.

## 6.2  Basic Definitions

### 6.2.1  Inner-product spaces

We recall that a real vector space [1] is a set, $H$, on which an addition and a scalar product are defined, namely $(h, h') \in H \times H \mapsto h + h' \in H$ and $(\lambda, h) \in \mathbb{R} \times H \mapsto \lambda h \in H$, and we assume that the reader is familiar with the theory of finite-dimensional

---

[1] All vector spaces in these notes will be real, and will therefore only be referred as vector spaces.

spaces.

An inner product on a vector space $H$ is a bilinear function, typically denoted $(\xi, \eta) \mapsto \langle \xi, \eta \rangle$ such that $\langle \xi, \xi \rangle \geq 0$ with $\langle \xi, \xi \rangle = 0$ if and only if $\xi = 0$. A vector space equipped with an inner product is called an inner-product space. We will often denote the inner product with a subscript referring to the space (e.g., $\langle \cdot, \cdot \rangle_H$). Given such a product, the function

$$\xi \mapsto \|\xi\|_H = \sqrt{\langle \xi, \xi \rangle_H}$$

is a norm, so that $H$ is also a normed space (but not all normed spaces are inner-product spaces)[2].

When a normed space is *complete* with respect to the topology induced by its norm, it is called a Banach space, or a Hilbert space when the norm is associated with an inner product. Completeness means that Cauchy sequences in this space always have a limit, i.e., if the sequence $(\xi_n)$ is such that, for any $\epsilon > 0$, there exists $n_0 > 0$ such that $\|\xi_n - \xi_m\|_H < \epsilon$ for all $n, m \geq n_0$, then there exists $\xi$ such that $\|\xi_n - \xi\|_H \to 0$. Completeness is a very natural property. It allows, for example, for the definition of integrals such as $\int h(t) dt$ as limits of Riemann sums for suitable functions $h: \mathbb{R} \to H$, leading (with more general notions of integrals) to proper definitions of expectations of $H$-valued random variables. Using a standard (abstract) construction, one can prove that any normed space (resp. inner-product) can be extended to a Banach (resp. Hilbert) space within which it is dense.

Note that finite-dimensional normed spaces are always complete.

### 6.2.2   Feature spaces and kernels

Now, consider an input set, say $\mathcal{R}$, and a mapping $h$ from $\mathcal{R}$ to $H$, where $H$ is an inner product space. For us, $\mathcal{R}$ is the set over which the original input data is observed, typically $\mathbb{R}^d$, and $H$ is the feature space. One can define the function $K_h: \mathcal{R} \times \mathcal{R} \to \mathbb{R}$ by

$$K_h(x, y) = \langle h(x), h(y) \rangle_H.$$

The function $K_h$ satisfies the following two properties.

[K1]  $K_h$ is symmetric, namely $K_h(x, y) = K_h(y, x)$ for all $x$ and $y$ in $\mathcal{R}$.

[K2]  For any $n > 0$, for any choice of scalars $\lambda_1, \ldots, \lambda_n \in \mathbb{R}$ and any $x_1, \ldots, x_n \in \mathcal{R}$, one has

$$\sum_{i,j=1}^{n} \lambda_i \lambda_j K_h(x_i, x_j) \geq 0. \tag{6.1}$$

---

[2]Note that we are using double bars for the norm in $H$, which, in most applications, is infinite dimensional

The first property is obvious, and the second one results from the fact that one can write

$$\sum_{i,j=1}^{n} \lambda_i \lambda_j K_h(x_i, x_j) = \sum_{i,j=1}^{n} \lambda_i \lambda_j \langle h(x_i), h(x_j)\rangle_H = \left\| \sum_{i=1}^{n} \lambda_i h(x_i) \right\|_H^2 \geq 0. \qquad (6.2)$$

This leads us to the following definition.

**Definition 6.1** *A function $K : \mathcal{R} \times \mathcal{R} \mapsto \mathbb{R}$ satisfying properties [K1] and [K2] is called a positive kernel.*

*One says that the kernel is positive definite if the sum in* (6.1) *cannot vanish unless (i) $\lambda_1 = \cdots = \lambda_n = 0$ or (ii) $x_i = x_j$ for some $i \neq j$.*

An equivalent definition of positive kernels can be given using kernel matrices, for which we introduce a notation.

**Definition 6.2** *If $K : \mathcal{R} \times \mathcal{R} \mapsto \mathbb{R}$ is given, we define, for every $x_1, \ldots, x_n \in \mathcal{R}$, the kernel matrix $\mathcal{K}_K(x_1, \ldots, x_n)$ with entries $K(x_i, x_j)$, for $i, j = 1, \ldots, n$. (If $K$ is understood from the context, we will simply write $\mathcal{K}(x_1, \ldots, x_n)$ instead of $\mathcal{K}_K(x_1, \ldots, x_n)$.)*

Given this notation, it is clear that $K$ is a positive kernel if and only if for all $x_1, \ldots, x_n \in \mathcal{R}$, the matrix $\mathcal{K}_K(x_1, \ldots, x_n)$ is symmetric, positive semidefinite. It is a positive definite kernel if $\mathcal{K}_K(x_1, \ldots, x_n)$ is positive definite as soon as all $x_j$'s are distinct. This latter condition is obviously needed since, if $x_i = x_j$, the $i$th and $j$th columns of $\mathcal{K}$ coincide and this matrix cannot be full-rank.

**Remark 6.3** It is important to point out that $K$ being a positive kernel *does not require* that $K(x, y) \geq 0$ for all $x, y \in \mathcal{R}$ (see examples in the next section). However, it does imply that $K(x, x) \geq 0$ for all $x \in \mathcal{R}$, since diagonal elements of positive semi-definite matrices are non-negative. ♦

The function $K_h$ defined above is therefore always a positive kernel, but not always positive definite, as seen below. We will also see later that the converse statement is true: any positive kernel $K : \mathcal{R} \times \mathcal{R} \mapsto \mathbb{R}$ can be expressed as $K_h$ for some feature function $h$ between $\mathcal{R}$ and some feature space $H$.

Given a feature function $h : \mathcal{R} \to H$, we will denote by $V_h = \text{span}(h(x), x \in \mathcal{R})$ the vector space generated by the features, which, by definition, is the space of all linear combinations

$$\xi = \sum_{i=1}^{n} \lambda_i h(x_i)$$

with $\lambda_1, \ldots, \lambda_m \in \mathbb{R}$, $x_1, \ldots, x_n \in \mathcal{R}$ and $n \geq 0$ (by convention, $\xi = 0$ if $n = 0$). Then $K_h$ is positive definite if and only if any family $(h(x_1), \ldots, h(x_n))$ with distinct $x_i$'s is linearly independent. This is a direct consequence of (6.2).

$$\sum_{i,j=1}^{n} \lambda_i \lambda_j K_h(x_i, x_j) = \left\| \sum_{i=1}^{n} \lambda_i h(x_i) \right\|_H^2 .$$

This implies in particular that positive-definite kernels over infinite input spaces $\mathcal{R}$ can only be associated to infinite-dimensional spaces $H$, since $V_h \subset H$.

## 6.3  First examples

### 6.3.1  Inner product

Clearly, if $\mathcal{R}$ is an inner product space, it has an associated reproducing kernel, defined by

$$K(x, y) = \langle x, y \rangle_{\mathcal{R}} .$$

This kernel is equal to $K_h$ with $H = \mathcal{R}$ and $h = \mathrm{id}$ (the identity mapping). In particular $K(x, y) = x^T y$ is a positive kernel if $\mathcal{R} = \mathbb{R}^d$. This kernel can obviously take positive and negative values.

Notice that this kernel is not positive definite, because the rank of $\mathcal{K}(x_1, \ldots, x_n)$ is equal to the dimension of $\mathrm{span}(x_1, \ldots, x_n)$, which can be less than $n$ even when the $x_i$'s are distinct.

### 6.3.2  Polynomial Kernels

Consider $\mathcal{R} = \mathbb{R}^d$ and define

$$h(x) = (x^{(i_1)} \ldots x^{(i_k)}, 1 \leq i_1, \ldots, i_k \leq d),$$

which contains all products of degree $k$ formed from variables $x^{(1)}, \ldots, x^{(d)}$, i.e., all monomials of degree $k$ in $x$. This function takes its values in the space $H = \mathbb{R}^{N_k}$, where $N_k = d^k$. Using, in $H$, the inner product $\langle \xi, \eta \rangle_H = \xi^T \eta$, we have

$$
\begin{aligned}
K_h(x, y) &= \sum_{1 \leq i_1, \ldots, i_k \leq d} (x^{(i_1)} y^{(i_1)}) \cdots (x^{(i_k)} y^{(i_k)}) \\
&= (x^T y)^k .
\end{aligned}
$$

This provides the homogeneous polynomial kernel of order $k$.

If one now takes all monomials of order less than or equal to $k$, i.e.,

$$h(x) = (x^{(i_1)} \dots x^{(i_l)}, 1 \le i_1, \dots, i_l \le d, 0 \le l \le k),$$

which now takes values in a space of dimension $1 + d + \cdots + d^k$, the corresponding kernel is

$$K_h(x, y) = 1 + (x^T y) + \cdots + (x^T y)^k = \frac{(x^T y)^{k+1} - 1}{x^T y - 1}.$$

This provides a polynomial kernel of order $k$. It is important to notice here that, even though the dimension of the feature space increases exponentially in $k$, so that the computation of the feature function rapidly becomes intractable, the computation of the kernel itself remains a relatively mild operation.

One can make variations on this construction. For example, choosing any family $c_0, c_1, \dots, c_k$ of positive numbers, one can take

$$h(x) = (c_l x^{(i_1)} \dots x^{(i_l)}, 1 \le i_1, \dots, i_l \le d, 0 \le l \le k)$$

yielding

$$K_h(x, y) = c_0^2 + c_1^2 (x^T y) + \cdots + c_k^2 (x^T y)^k.$$

Taking $c_l = \binom{k}{l}^{1/2} \alpha^l$ for some $\alpha > 0$, we get another form of polynomial kernel, namely,

$$K_h(x, y) = (1 + \alpha^2 x^T y)^k.$$

### 6.3.3  Functional Features

We now consider an example in which $H$ is infinite dimensional. Let $\mathcal{R} = \mathbb{R}^d$. We assume that a function $s : \mathbb{R}^d \to \mathbb{R}$ is chosen, such that $s$ is both (absolutely) integrable and square integrable. We also fix a scaling parameter $\rho > 0$. Associate to $x \in \mathbb{R}^d$ the function

$$\xi_x : y \mapsto s((y - x)/\rho),$$

which is also square integrable (as a function of $y$). We define the feature function $h : x \mapsto \xi_x$ from $\mathbb{R}^d$ to $H = L^2(\mathbb{R}^d)$, the space of square integrable functions on $\mathbb{R}^d$ with inner product

$$\langle \xi, \eta \rangle_H = \int_{\mathbb{R}^d} \xi(z) \eta(z) dz.$$

The resulting kernel is

$$K_h(x, y) = \int_{\mathbb{R}^d} s(z/\rho - x) s(z/\rho - y) \, dz = \rho^d \int_{\mathbb{R}^d} s(z) s(z - (y - x)/\rho) \, dz.$$

Note that $K_h(x,y)$ is "translation-invariant," which means that it only depends on $x - y$. It takes the form $K_h(x,y) = \rho^d \Gamma((y-x)/\rho)$ where

$$\Gamma(u) = \int_{\mathbb{R}^d} s(z)s(z-u)\,dz.$$

is the convolution[3] of $s$ with $\tilde{s} : z \mapsto s(-z)$.

Let $\sigma$ be the Fourier transform of $s$, i.e.,

$$\sigma(\omega) = \int_{\mathbb{R}^d} e^{-2i\pi\omega^T u} s(u)\,du.$$

Because $s$ is real-valued, we have $\sigma(-\omega) = \bar{\sigma}(\omega)$, the complex conjugate of $\sigma$. Moreover, $\bar{\sigma}$ is also the Fourier transform of $\tilde{s}$. Using the fact that the Fourier transform of the convolution of two functions is the product of their Fourier transforms, we see that the Fourier transform of $\Gamma = s * \tilde{s}$ is equal to $|\sigma|^2$. Applying the inverse transform, we find

$$\Gamma(u) = \int_{\mathbb{R}^d} e^{2i\pi\omega^T u} |\sigma(\omega)|^2\,d\omega = \int_{\mathbb{R}^d} e^{-2i\pi\omega^T u} |\sigma(-\omega)|^2\,d\omega.$$

This form is (almost) characteristic of translation-invariant kernels.

Let us consider a few examples of kernels that can be obtained in this way.

(1) Take $d = 1$ and let $s$ be the indicator function of the interval $[-\frac{1}{2}, \frac{1}{2}]$. Then, one finds

$$\Gamma(t) = \max(1 - |t|, 0).$$

In this case, the space $V_h$ is the space of all functions expressed as finite sums

$$z \mapsto \sum_{j=1}^{n} \lambda_j \mathbf{1}_{[x_j - \rho/2, x_j + \rho/2]}(z),$$

and therefore is a space of compactly-supported piecewise constant functions. Such a function computed with distinct $x_j$'s cannot vanish everywhere unless all $\lambda_j$'s vanish, so that $K_h$ is positive definite. Indeed, let

$$f(z) = \sum_{j=1}^{n} \lambda_j \mathbf{1}_{[x_j - \rho/2, x_j + \rho/2]}(z)$$

and assume without loss of generality that $x_1 < x_2 < \cdots < x_n$ and let $x_{n+1} = \infty$. Let $i_0$ be the smallest index $j$ such that $\lambda_j \neq 0$, assuming that such an index exists. Then $f(z) = \lambda_{i_0} > 0$ for all $z \in [x_{i_0} - \rho/2, x_{i_0+1} - \rho/2)$ which is a non-empty interval. So, if $f$ vanishes almost everywhere, we must have $\lambda_j = 0$ for all $j = 1, \ldots, n$.

---

[3]The convolution between two absolutely integrable functions $f$ and $g$ is defined by $f * g(u) = \int_{\mathbb{R}^d} f(z)g(u-z)\,dz$

(2) Still with $d = 1$, let $s(z) = e^{-|z|}$. Then, for $t > 0$,

$$\Gamma(t) = \int_{-\infty}^{\infty} e^{-|z|} e^{-|z-t|} \, dz$$

$$= \int_{-\infty}^{0} e^z e^{z-t} \, dz + \int_{0}^{t} e^{-z} e^{z-t} \, dz + \int_{t}^{\infty} e^{-z} e^{-z+t} \, dz$$

$$= \frac{e^{-t}}{2} + t e^{-t} + \frac{e^{-t}}{2}$$

$$= (1 + t) e^{-t}$$

Using the fact that $\Gamma(-t) = \Gamma(t)$ (make the change of variable $z \to -z$ in the integral), we get

$$\Gamma(t) = (1 + |t|) e^{-|t|}.$$

for all $t$. This shows that

$$K(x, y) = (1 + |x - y|) e^{-|x-y|}$$

is a positive kernel on $\mathbb{R}^d$.

(3) Take $s(z) = e^{-|z|^2/2}$, $z \in \mathbb{R}^d$. Then

$$\Gamma(u) = \int_{\mathbb{R}^d} e^{-\frac{|z|^2 + |u-z|^2}{2}} \, dz = e^{-\frac{|u|^2}{4}} \int_{\mathbb{R}^d} e^{-|z-u/2|^2} \, dz$$

$$= (4\pi)^{d/2} e^{-\frac{|u|^2}{4}}.$$

This provides a special case of Gaussian kernel.

### 6.3.4 General construction theorems

**Translation invariance**

As introduced above, a kernel $K$ is translation invariant if it takes the form $K(x, y) = \Gamma(x-y)$ for some continuous function $\Gamma$ defined on $\mathbb{R}^d$. Bochner's theorem [33] states that such a $K$ is a positive kernel if and only if $\Gamma$ is the Fourier transform of a positive measure, namely,

$$\Gamma(x) = \int_{\mathbb{R}^d} e^{-2i\pi\langle x, \omega \rangle} d\mu(\omega)$$

where $\mu$ is a positive and symmetric (invariant by sign change) measure on $\mathbb{R}^d$. For example one can take $d\mu(\omega) = \nu(\omega) d\omega$, where $\nu$ is a integrable, positive and even function.

This theorem provides an at least numerical, and sometimes analytical, method for constructing kernels. The previous section exhibited a special case of translation-invariant kernel for which $\nu = |\sigma|^2$.

**Radial kernels**

A radial kernel takes the form $K(x,y) = \gamma(|x-y|^2)$, for some continuous function $\gamma$ defined on $[0,+\infty)$. Shoenberg's theorem [173] states that, if this function $\gamma$ is universally valid, i.e., $K$ is a kernel for all dimensions $d$, then, it must take the form

$$\gamma(t) = \int_0^\infty e^{-\lambda t} d\mu(\lambda)$$

for some positive finite measure $\mu$ on $[0,+\infty)$.

For example, when $\mu$ is a Dirac measure, i.e., $\mu = \delta_{(2a)^{-1}}$ for some $a > 0$, then $K(x,y) = \exp(-|x-y|^2/2a)$, which is the Gaussian kernel. Taking $d\mu = e^{-a\lambda}d\lambda$ yields $\gamma(t) = 1/(t+a)$, and $d\mu = \lambda e^{-a\lambda}d\lambda$ yields $\gamma(t) = 1/(a+t)^2$.

There is also, in Schoenberg [173], a characterization of radial kernels for a fixed dimension $d$. Such kernels must take the form

$$\gamma(t) = \int_0^{+\infty} \Omega_d(t\lambda)d\mu(\lambda)$$

with $\Omega_d(t) = \Gamma(d/2)(2/t)^{(d-2)/2}J_{(d-2)/2}(t)$ where $J_{(d-2)/2}$ is Bessel's function of the first kind.

### 6.3.5   Operations on kernels

Kernels can be combined in several ways as described in the next proposition.

**Proposition 6.4** *Let $K_1 : \mathcal{R} \times \mathcal{R} \to \mathbb{R}$ and $K_2 : \mathcal{R} \times \mathcal{R} \to \mathbb{R}$ be positive kernels. Then the following assertions hold.*

(i) *If $\lambda_1, \lambda_2 > 0$, $\lambda_1 K_1 + \lambda_2 K_2$ is a positive kernel. It is positive definite as soon as either $K_1$ or $K_2$ is positive definite.*

(ii) *For any function $f : \mathcal{R}' \to \mathcal{R}$, $K_1'(x',y') \overset{\Delta}{=} K_1(f(x'), f(y'))$ is a positive kernel. It is positive definite as soon as $K_1$ is positive definite and $f$ is one-to-one.*

(iii) *$K(x,y) = K_1(x,y)K_2(x,y)$ is a positive kernel. It is positive definite as soon as $K_1$ and $K_2$ are positive definite.*

(iv) *Let $K_1$ and $K_2$ be translation-invariant with $\mathcal{R} = \mathbb{R}^d$, taking the form $K_i(x,y) = \Gamma_i(x-y)$, where $\Gamma_i$ is continuous ( $i = 1,2$). Assume that one of the two functions $\Gamma_1, \Gamma_2$ is integrable on $\mathbb{R}^d$. Then*

$$K(x,y) = \int_{\mathbb{R}^d} K_1(x,z)K_2(z,y)dz$$

*is also a positive kernel.*

PROOF Point (i) is obvious. Point (ii) is almost as simple, because, for any $\lambda_1, \ldots, \lambda_n \in \mathbb{R}$ and $x'_1, \ldots, x'_n \in \mathcal{R}'$,

$$\sum_{i,j=1}^{n} \lambda_i \lambda_j K'_1(x'_i, x'_j) = \sum_{i,j=1}^{n} \lambda_i \lambda_j K_1(f(x'_i), f(x'_j)) \geq 0.$$

If $K_1$ is positive definite, then the latter sum can only vanish if all $\lambda_i$ are zero, or some of the points in $(f(x'_1), \ldots, f(x'_n))$ coincide. If, in addition, $f$ is one-to-one, then this is equivalent to all $\lambda_i$ are zero, or some of the points in $(x'_1, \ldots, x'_n)$ coincide, so that $K'_1$ is positive definite.

To prove point (iii), take $x_1, \ldots, x_N \in \mathbb{R}^d$ and form the matrices $\mathcal{K}_i = \mathcal{K}_i(x_1, \ldots, x_N)$, $i = 1, 2$, which are, by assumption positive semi-definite. The matrix $\mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$ is the element-wise (or Hadamard) product of $\mathcal{K}_1$ and $\mathcal{K}_2$, and the conclusion follows from the linear algebra result stating that the Hadamard product of two positive semi-definite (resp. positive definite) matrices $A = (a(i, j), 1 \leq i, j \leq N)$ and $B = (b(i, j), 1 \leq i, j \leq N)$ is positive semi-definite (resp. positive definite). This is proved by diagonalizing, say, $A$ in an orthonormal basis $u_1, \ldots, u_N$, with eigenvalues $\lambda_1, \ldots, \lambda_N$ and writing

$$\sum_{i,j=1}^{N} \alpha^{(i)} a(i, j) b(i, j) \alpha^{(j)} = \sum_{i,j,k=1}^{N} \alpha^{(i)} u_i^{(k)} u_j^{(k)} \lambda_k b(i, j) \alpha^{(j)}$$

$$= \sum_{k=1}^{N} \lambda_k \sum_{i,j=1}^{N} (\alpha^{(i)} u_i^{(k)})(\alpha^{(j)} u_j^{(k)})) b(i, j) \geq 0$$

If $B$ is positive definite, then the sum above can be zero only if, for each $k$, either $\lambda_k = 0$ or $\alpha^{(i)} u_i^{(k)} = 0$ for all $i$. If $A$ is also positive definite, then the only possibility is $\alpha^{(i)} u_i^{(k)} = 0$ for all $i$ and $k$, which implies $\alpha^{(i)} = 0$ for all $i$ since $u_i \neq 0$.

To prove point (iv) [4], we first note that a translation invariant kernel $K'(x, y) = \Gamma'(x - y)$ is always bounded. Indeed, the matrix $\mathcal{K}'(x, 0)$ is positive semi-definite, with determinant $\Gamma'(0)^2 - \Gamma'(x)^2 > 0$, showing that $|\Gamma'(x)| < \Gamma'(0)$. This shows that the integral defining $K(x, y)$ converges as soon as one of the two functions $\Gamma_1$ or $\Gamma_2$ is integrable. Moreover, we have $K(x, y) = \Gamma(x - y)$ with

$$\Gamma(x) = \int_{\mathbb{R}^d} \Gamma_1(x - z) \Gamma_2(z) \, dz = \int_{\mathbb{R}^d} \Gamma_1(x - u) \Gamma_2(u - y) \, du$$

Using the fact that both $\Gamma_1$ and $\Gamma_2$ are even, and making the change of variable $z \mapsto -z$, one easily shows that $\Gamma(x) = \Gamma(-x)$, which implies that $K$ is symmetric.

---

[4]This part of the proof uses some measure theory.

We proceed with the assumption that $\Gamma_2$ is integrable and use Bochner's theorem to write

$$\Gamma_1(y) = \int_{\mathbb{R}^d} e^{-i\xi^T y} d\mu_1(\xi)$$

for some positive finite measure $\mu_1$. Then

$$\Gamma(x) = \int_{\mathbb{R}^d} \left( \int_{\mathbb{R}^d} e^{-2i\pi\xi^T(x-z)} d\mu_1(\xi) \right) \Gamma_2(z) \, dz$$

$$= \int_{\mathbb{R}^d} e^{-2i\pi\xi^T x} \left( \int_{\mathbb{R}^d} e^{2i\pi\xi^T z} \Gamma_2(z) \, dz \right) d\mu_1(\xi)$$

The shift in the order of the variables $\xi$ and $z$ uses Fubini's theorem. The function

$$\psi(\xi) = \int_{\mathbb{R}^d} e^{2i\pi\xi^T z} \Gamma_2(z) \, dz$$

is the inverse Fourier transform of $\Gamma_2$. Because $\Gamma_2$ is bounded and integrable, it is also square integrable, which implies that its inverse Fourier transform is also a square integrable function. Since Bochner's theorem implies that $\Gamma_2$ is the Fourier transform of a positive measure $\mu_2$, we find, using the injectivity of the Fourier transform, that $\psi$ is non-negative. So $\Gamma$ is the Fourier transform of the finite positive measure $\psi d\mu_1$, which implies that $K$ is a positive kernel. ∎

Point (iv) can be related to the following discrete statement on symmetric matrices: assume that $A$ and $B$ are positive semi-definite and that they commute, so that $AB = BA$: then $AB$ is positive semi-definite (see **??**). In the case of kernels, one may consider the symmetric linear operators $\mathbb{K}_i : f \mapsto \int_{\mathbb{R}^d} K_i(\cdot, y) f(y) dy$ which maps the space of square integrable functions into itself. Then $\mathbb{K}_1$ and $\mathbb{K}_2$ commute and $\mathbb{K} = \mathbb{K}_1 \mathbb{K}_2$.

### 6.3.6 Canonical Feature Spaces

Let $K$ be a positive kernel on a set $\mathcal{R}$. The following construction, which is fundamental, shows that $K$ can always be associated with a feature function $h$ taking values in a suitably chosen inner-product space $H$.

Associate to each $x \in \mathcal{R}$ the function $\xi_x : y \mapsto K(y, x)$ (we will also write $\xi_x = K(\cdot, x)$), and let $H_K = \mathrm{span}(\xi_x, x \in \mathcal{R})$, a subspace of the vector space of all functions from $\mathcal{R}$ to $\mathbb{R}$. Define the feature function $h : x \mapsto \xi_x$ from $\mathcal{R}$ to $H_K$. There is a unique inner product on $H_K$ such that $K = K_h$. Indeed, by definition, this requires

$$\langle K(\cdot, x), K(\cdot, y) \rangle_{H_K} = K(x, y). \tag{6.3}$$

Moreover, by linearity, for any $\xi = \sum_{i=1}^{n} \lambda_i K(\cdot, x_i)$ and $\eta = \sum_{i=1}^{m} \mu_i K(\cdot, y_i)$, one needs

$$\langle \xi, \eta \rangle_{H_K} = \sum_{i=1}^{n} \sum_{j=1}^{m} \lambda_i \mu_j K(x_i, y_j),$$

so that the inner product is uniquely specified on $H_K$. To make sure that this inner-product is well defined, we must check that there is no ambiguity, in the sense that, if $\xi$ has an alternative decomposition $\xi = \sum_{i=1}^{n'} \lambda_i' K(\cdot, x_i')$, then, the value of $\langle \xi, \eta \rangle_{H_K}$ remains unchanged. But this is clear, because one can also write

$$\langle \xi, \eta \rangle_{H_K} = \sum_{j=1}^{m} \mu_j \xi(y_j),$$

which only depends on $\xi$ and not on its decomposition. The linearity of the product with respect to $\xi$ is also clear from this expression, and the bilinearity by symmetry.

The Schwartz inequality implies that

$$|\langle \xi, \eta \rangle_{H_K}| \le \|\xi\|_{H_K} \|\eta\|_{H_K}$$

From which we deduce that $\|\xi\|_{H_K} = 0$ implies that $\langle \xi, \eta \rangle_{H_K} = 0$ for $\eta \in H_K$. Since $\langle \xi, K(\cdot, y) \rangle_{H_K} = \xi(y)$ for all $y$, this also implies that $\xi = 0$, completing the proof that $H_K$ is an inner-product space.

Equation (6.3) is the "reproducing property" of the kernel for the inner-product on $H_K$. In functional analysis, the completion, $\hat{H}_K$, of $H_K$ for the topology associated to its norm is then a Hilbert space, and is referred to as a "reproducing kernel Hilbert space," or RKHS.

More generally, an inner-product space $H$ of functions $h : \mathcal{R} \to \mathbb{R}$ is a reproducing kernel Hilbert space if $H$ is a complete space (which makes it Hilbert) and there exists a positive kernel $K$ such that,

[RKHS1]  For all $x \in \mathcal{R}$, $K(\cdot, x)$ belongs to $H$,

[RKHS2]  For all $h \in H$ and $x \in \mathcal{R}$,

$$\langle h, K(\cdot, x) \rangle_H = h(x).$$

Returning to the example of functional features in section 6.3.3, we have two different representations of the kernel in feature space, namely in $H = L^2(\mathbb{R}^d)$, or in $H_K$, with a different inner product. There is not a contradiction, and simply shows that the representation of a positive kernel in terms of a feature function is not unique.

## 6.4   Projection on a finite-dimensional subspace

If $H$ is an inner-product space and $V$ is a subspace of $H$, one defines the orthogonal projection of an element $\xi \in H$ on $V$ as its closest point in $V$, that is, the element $\eta^*$ of $V$ minimizing the function $F : \eta \mapsto \|\eta - \xi\|_H^2$ over all $\eta \in V$. This closest point does not always exist, but it does in the special case in which $V$ is finite dimensional (or, more generally, when $V$ is a closed subspace of $H$; see Yosida [205]). We state, without proof, some of the properties of this operation.

Assuming that $V$ is closed, this minimizer is unique and will be denoted $\eta^* = \pi_V(\xi)$. Moreover, $\pi_V$ is a linear transformation from $H$ to $V$, and $\eta^*$ is characterized by the properties

$$\begin{cases} \eta^* \in V \\ \xi - \eta^* \perp V, \end{cases}$$

the last condition meaning that $\langle \xi - \eta^*, \eta \rangle_H = 0$ for all $\eta \in V$.

Because $\|\xi\|_H^2 = \|\pi_V(\xi)\|_H^2 + \|\xi - \pi_V(\xi)\|_H^2$, one always has $\|\pi_V(\xi)\|_H \le \|\xi\|_H$, with inequality if and only if $\pi_V(\xi) = \xi$, i.e., if and only if $\xi \in V$.

If $V$ is finite-dimensional and $\eta_1, \ldots, \eta_n$ is a basis of $V$, then $\pi_V(\xi)$ is given by

$$\pi_V(\xi) = \sum_{i=1}^n \alpha^{(i)} \eta_i$$

with $\alpha$ (considered as a column vector in $\mathbb{R}^n$) given by

$$\alpha = \mathrm{Gram}(\eta_1, \ldots, \eta_n)^{-1} \lambda,$$

where $\lambda \in \mathbb{R}^n$ is the vector with coordinates $\lambda^{(i)} = \langle \xi, \eta_i \rangle_H$, $i = 1, \ldots, n$. The Gram matrix of $\eta_1, \ldots, \eta_n$, denoted $\mathrm{Gram}(\eta_1, \ldots, \eta_n)$, is the $n$ by $n$ matrix with entries $\langle \eta_i, \eta_j \rangle_H$ for $i, j = 1, \ldots, n$.

If $A$ is a subset of $H$, the set $A^\perp$ consists of all vectors perpendicular to $A$, namely

$$A^\perp = \left\{ h \in H : \langle h, \tilde{h} \rangle_H = 0 \text{ for all } \tilde{h} \in A \right\}.$$

If $V$ is a finite-dimensional (or, more generally, closed) subspace of $H$, then any point in $h$ is decomposed as $h = \pi_V(h) + h - \pi_V(h)$ with $h - \pi_V(h) \in V^\perp$. This shows that $\pi_{V^\perp}$ is well defined and equal to $\mathrm{id}_H - \pi_V$.

Orthogonal projections can be applied to function interpolation in an RKHS. Indeed, assuming that $H$ is an RKHS, as described at the end of the previous section, with a positive-definite kernel. Given distinct points $x_1, \ldots, x_N \in \mathcal{R}$ and values

$\alpha_1, \ldots, \alpha_N \in \mathbb{R}$, the interpolation problem consists in finding $h \in H$ with minimal norm satisfying $h(x_k) = \alpha_k$, $k = 1, \ldots, N$. Consider the finite dimensional space

$$V = \mathrm{span}\{K(\cdot, x_k), k = 1, \ldots N\}.$$

Then there exists an element $h_0 \in V$ that satisfies the constraints. Indeed, looking for $h_0$ in the form

$$h_0(x) = \sum_{l=1}^{N} K(x, x_l)\lambda_l$$

one has

$$h_0(x_k) = \sum_{l=1}^{N} K(x_k, x_l)\lambda_l$$

so that

$$\begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_N \end{pmatrix} = \mathcal{K}(x_1, \ldots, x_N)^{-1} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix}$$

Any other function $h$ satisfying the constraints satisfies $h(x_k) - h_0(x_k) = 0$, which, using RKHS2, is equivalent to $\langle h - h_0, K(\cdot, x_k)\rangle_H = 0$, i.e., to $h - h_0 \in V^\perp$. This shows that $h_0 = \pi_V(h)$, so that $\|h\|_H \geq \|h_0\|_H$ and $h_0$ provides the optimal interpolation. We summarize this in the proposition:

**Proposition 6.5** *Let H is an RKHS with a positive-definite kernel. Let $x_1, \ldots, x_N \in \mathcal{R}$ be distinct points and $\alpha_1, \ldots, \alpha_N \in \mathbb{R}$. Then the function $h \in H$ with minimal norm satisfying $h(x_k) = \alpha_k$, $k = 1, \ldots, N$ takes the form*

$$h(x_k) = \sum_{l=1}^{N} K(x_k, x_l)\lambda_l \tag{6.4a}$$

*with*

$$\begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_N \end{pmatrix} = \mathcal{K}(x_1, \ldots, x_N)^{-1} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix}. \tag{6.4b}$$

A variation of this problem replaces the constraint by a penalty that complete the minimization associated with the orthogonal projection, namely, minimizing (in $h \in H$)

$$\|h\|_H^2 + \sigma^2 \sum_{k=1}^{N} |h(x_k) - \alpha_k|^2.$$

Letting $h_0 = \pi_V(h)$, so that $h_0(x_k) = h(x_k)$ for all $k$, this expression can be rewritten as

$$\|h_0\|_H^2 + \|h - h_0\|_H^2 + \sigma^2 \sum_{k=1}^{N} |h_0(x_k) - \alpha_k|^2.$$

This shows that the optimal $h$ must coincide with its projection on $V$, and therefore belong to that subspace. Looking for $h$ in the form

$$h(\cdot) = \sum_{l=1}^{N} K(\cdot, x_l)\lambda_l,$$

the objective function is rewritten as

$$\sum_{k,l=1}^{N} K(x_k, x_l)\lambda_k\lambda_l + \sigma^2 \sum_{k=1}^{N} \left| \sum_{l=1}^{N} K(x_k, x_l)\lambda_l - \alpha_k \right|^2,$$

which, in vector notation gives, writing $\lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_N \end{pmatrix}$ and $\alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix}$,

$$\lambda^T \mathcal{K}(x_1, \ldots, x_N)\lambda + \sigma^2 (\mathcal{K}(x_1, \ldots, x_N)\lambda - \alpha)^T (\mathcal{K}(x_1, \ldots, x_N)\lambda - \alpha).$$

The differential of this expression in $\lambda$ is

$$\mathcal{K}(x_1, \ldots, x_N)\lambda + 2\sigma^2 \mathcal{K}(x_1, \ldots, x_N)(\mathcal{K}(x_1, \ldots, x_N)\lambda - \alpha).$$

Assuming that $x_1, \ldots, x_N$ are distinct, this vanishes if and only if

$$\lambda = (\mathcal{K}(x_1, \ldots, x_N) + (1/\sigma^2)\mathrm{Id}_{\mathbb{R}^N})^{-1}\alpha.$$

We have just proved the proposition:

**Proposition 6.6** *Let $H$ is an RKHS with a positive-definite kernel. Let $x_1, \ldots, x_N \in \mathcal{R}$ be distinct points and $\alpha_1, \ldots, \alpha_N \in \mathbb{R}$. Then the unique minimizer of*

$$h \mapsto \|h\|_H^2 + \sigma^2 \sum_{k=1}^{N} |h(x_k) - \alpha_k|^2$$

*on $H$ is given by*

$$h(x_k) = \sum_{l=1}^{N} K(x_k, x_l)\lambda_l \tag{6.5a}$$

*with*

$$\begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_N \end{pmatrix} = (\mathcal{K}(x_1, \ldots, x_N) + (1/\sigma^2)\mathrm{Id}_{\mathbb{R}^N})^{-1} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix}. \tag{6.5b}$$

# Chapter 7

# Linear Models for Regression

In regression, *linear models* refer to situations in which one tries to predict the dependent variable $Y \in \mathcal{R}_Y = \mathbb{R}^q$ by a function $\hat{f}(X)$ of the dependent variable $X \in \mathcal{R}_X$, where $\hat{f}$ is optimized over a *linear space* $\mathcal{F}$. The most common situation is the "standard linear model," for which $\mathcal{R}_X = \mathbb{R}^d$ and

$$\mathcal{F} = \{f(x) = a_0 + b^T x : a_0 \in \mathbb{R}^q, b \in \mathcal{M}_{d,q}(\mathbb{R})\}. \tag{7.1}$$

More generally, with $q = 1$, given a mapping $h : \mathcal{R} \to H$, where $H$ is an inner-product space, one can take:

$$\mathcal{F} = \{f(x) = a_0 + \langle b, h(x) \rangle_H : a_0 \in \mathbb{R}, b \in H\}. \tag{7.2}$$

Note that $h$ can be nonlinear, and $\mathcal{F}$ can be infinite dimensional. Such sets corresponds to linear models using feature functions, and will be addressed using kernel methods in this chapter.

Note also that, even if the model is linear, the associated training algorithms can be nonlinear, and we will review in fact several situations in which solving the estimation problem requires nonlinear optimization methods.

## 7.1 Least-Square Regression

### 7.1.1 Notation and Basic Estimator

We denote by $Y$ and $X$ the dependent and independent variables of the regression problem. We will assume that $Y$ takes values in $\mathbb{R}^q$ and that $X$ takes values in a set $\mathcal{R}_X$, which will, by default, be equal to $\mathbb{R}^d$, except when discussing kernel methods, for which this set can be arbitrary (provided that there is a mapping $h$ from $\mathcal{R}_X$ to an inner product space $H$ with an easily computable kernel).

Least-square regression uses the risk function $r(y, y') = |y - y'|^2$. The prediction error is then $R(f) = E(|Y - f(X)|^2)$ for any predictor $f$ and the Bayes predictor is the conditional expectation $x \mapsto E(Y \mid X = x)$ (see item Example 1. in section 5.3). We also start with the standard setting where $\mathcal{R}_X = \mathbb{R}^d$ and $\mathcal{F}$ given by (7.1).

We will use the following notation, which sometimes simplifies the computation. If $x \in \mathbb{R}^d$, we let $\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$, which belongs to $\mathbb{R}^{d+1}$. The linear predictor $f(x) = a_0 + b^T x$ with $a_0 \in \mathbb{R}^q, b \in \mathcal{M}_{d,q}(\mathbb{R})$ can then be written as $f(x) = \beta^T \tilde{x}$ with $\beta = \begin{pmatrix} a_0^T \\ b \end{pmatrix} \in \mathcal{M}_{d+1,q}(\mathbb{R})$.

In a model-based approach, the linear model is a Bayes predictor under the generative assumption that $Y = a_0 + b^T X + \epsilon$ where $\epsilon$ is a residual noise satisfying $E(\epsilon \mid X) = 0$, which is true, for example, when $\epsilon$ is centered and independent of $X$. If one further specifies the model so that $\epsilon$ is Gaussian, centered and independent of $X$, and one assumes that the distribution of $X$ does not depend on $a_0$ and $b$, then the maximum likelihood estimator of these parameters based on a training set $T = ((x_1, y_1), \ldots, (x_N, y_N))$ must minimize the "residual sum of squares:"

$$RSS(\beta) \triangleq N\hat{R}(f) \triangleq \sum_{k=1}^{N} |y_k - f(x_k)|^2 = \sum_{k=1}^{N} |y_k - \beta^T \tilde{x}_k|^2.$$

In other terms, the model-based approach is identical, under these (standard) assumptions, to empirical risk minimization (section 5.5), on which we now focus. (Recall that, even when using a model-based approach, one does not make assumptions on the true distribution of $X$ and $Y$; one rather treats the model as an approximation of these distributions, estimated by maximum likelihood, and uses the Bayes predictor for the estimated model.)

The computation of the optimal regression parameters is made easier by the introduction of the following matrices. Introduce the $N \times (d + 1)$ matrix $\mathcal{X}$ with rows $\tilde{x}_1^T, \ldots, \tilde{x}_N^T$ and the $N \times q$ matrix $\mathcal{Y}$ with rows $y_1^T, \ldots, y_N^T$, that is:

$$\mathcal{X} = \begin{pmatrix} 1 & x_1^{(1)} & \cdots & x_1^{(d)} \\ \vdots & \vdots & & \vdots \\ 1 & x_N^{(1)} & \cdots & x_N^{(d)} \end{pmatrix}, \quad \mathcal{Y} = \begin{pmatrix} y_1^{(1)} & \cdots & y_1^{(q)} \\ \vdots & & \vdots \\ y_N^{(1)} & \cdots & y_N^{(q)} \end{pmatrix}.$$

With this notation, we have

$$RSS(\beta) = |\mathcal{Y} - \mathcal{X}\beta|_2^2.$$

with $|A|_2^2 = \text{trace}(A^T A)$ for a rectangular matrix $A$. The solution of the problem is then provided by the following theorem.

**Theorem 7.1** *Assume that the matrix $\mathcal{X}$ has rank $d + 1$. Then the RSS is minimized for*

$$\hat{\beta} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathcal{Y}$$

PROOF We provide two possible proofs of this elementary problem. The first one is an optimization argument noting that $F(\beta) \overset{\Delta}{=} RSS(\beta)$ is a convex function defined on $\mathcal{M}_{d+1,q}(\mathbb{R})$ and with values in $\mathbb{R}$. Since $F$ is quadratic, we have, for any matrix $h \in \mathcal{M}_{d+1,q}(\mathbb{R})$,

$$dF(\beta)h = \partial_\epsilon F(\beta + \epsilon h)|_{\epsilon=0} = -2\mathrm{trace}(h^T \mathcal{X}^T (\mathcal{Y} - \mathcal{X}\beta))$$

and

$$dF(\beta) = 0 \Leftrightarrow \mathcal{X}^T (\mathcal{Y} - \mathcal{X}\beta) = 0 \Leftrightarrow \beta = \hat{\beta}.$$

One can alternatively proceed with a direct computation. We have

$$RSS(\beta) = |\mathcal{Y}|_2^2 - 2\mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{Y}) + \mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{X}\beta)$$
$$= |\mathcal{Y}|_2^2 - 2\mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{X}\hat{\beta}) + \mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{X}\beta).$$

Replacing $\beta$ by $\hat{\beta}$ and simplifying yields

$$RSS(\hat{\beta}) = |\mathcal{Y}|_2^2 - \mathrm{trace}(\hat{\beta}^T \mathcal{X}^T \mathcal{X}\hat{\beta})$$

It follows that

$$RSS(\beta) = RSS(\hat{\beta}) + \mathrm{trace}(\hat{\beta}^T \mathcal{X}^T \mathcal{X}\hat{\beta}) - 2\mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{X}\hat{\beta}) + \mathrm{trace}(\beta^T \mathcal{X}^T \mathcal{X}\beta)$$
$$= RSS(\hat{\beta}) + |\mathcal{X}(\hat{\beta} - \beta)|_2^2$$

so that the left-hand side is minimized at $\beta = \hat{\beta}$. ∎

**Remark 7.2** If $\mathcal{X}$ does not have rank $d + 1$, then optimal solutions exist, but they are not unique. By convexity, the solutions are exactly the vectors $\beta$ at which the gradient vanishes, i.e., those that satisfy $\mathcal{X}^T \mathcal{X}\beta = \mathcal{X}^T \mathcal{Y}$. The set of solutions can be obtained by introducing the SVD of $\mathcal{X}$ in the form $\mathcal{X} = UDV^T$ and letting $\gamma = V^T \beta$ and $\mathcal{Z} = U^T \mathcal{Y}$. Then

$$\mathcal{X}^T \mathcal{X}\beta = \mathcal{X}^T \mathcal{Y} \Leftrightarrow D^T D\gamma = D^T \mathcal{Z}.$$

Letting $d^{(1)}, \ldots, d^{(m)}$ denote the nonzero diagonal entries of $D$ (so that $m \leq d + 1$), we find $\gamma^{(i)} = z^{(i)}/d^{(i)}$ for $i \leq m$ (the other equalities being $0 = 0$). So, the $d + 1 - m$ last entries of $\gamma$ can be chosen arbitrarily (and $\beta = V\gamma$). ♦

An alternate representation of the solution use a two-step computation that estimates $b$ first, then $a_0$. Indeed, for fixed $\hat{b}$, the minimum of

$$\sum_{k=1}^{N} |y_k - a_0 - x_k^T \hat{b}|^2$$

is attained at $\hat{a}_0 = \bar{y} - \bar{x}^T \hat{b}$ with the usual definitions

$$\bar{y} = \frac{1}{N} \sum_{k=1}^{N} y_k \text{ and } \bar{x} = \frac{1}{N} \sum_{k=1}^{N} x_k.$$

This shows that $\hat{b}$ itself must be a minimizer of

$$\sum_{k=1}^{N} |y_k - \bar{y} - (x_k - \bar{x})^T b|^2.$$

Denote by $\mathcal{Y}_c$ and $\mathcal{X}_c$ the matrices

$$\mathcal{X}_c = \begin{pmatrix} x_1^{(1)} - \bar{x}^{(1)} & \cdots & x_1^{(d)} - \bar{x}^{(d)} \\ \vdots & & \vdots \\ x_N^{(1)} - \bar{x}^{(1)} & \cdots & x_N^{(d)} - \bar{x}^{(d)} \end{pmatrix}, \mathcal{Y}_c = \begin{pmatrix} y_1^{(1)} - \bar{y}^{(1)} & \cdots & y_1^{(q)} - \bar{y}^{(q)} \\ \vdots & & \vdots \\ y_N^{(1)} - \bar{y}^{(1)} & \cdots & y_N^{(q)} - \bar{y}^{(q)} \end{pmatrix}.$$

Then $\hat{b}$ must minimize $|\mathcal{Y}_c - \mathcal{X}_c b|^2$, yielding

$$\hat{b} = (\mathcal{X}_c^T \mathcal{X}_c)^{-1} \mathcal{X}_c^T \mathcal{Y}_c, \quad \hat{a}_0 = \bar{y} - \bar{x}^T \hat{b}.$$

The reader may want to double-check that this solution coincides with the one provided in theorem 7.1.

### 7.1.2  Limit behavior

The matrix

$$\hat{\Sigma}_{XX} = \frac{1}{N} \mathcal{X}_c^T \mathcal{X}_c = \frac{1}{N} \sum_{k=1}^{N} (x_k - \bar{x})(x_k - \bar{x})^T$$

is a sample estimate of the covariance matrix of $X$, that we will denote $\Sigma_{XX}$. Similarly, $\hat{\Sigma}_{XY} = \mathcal{X}_c^T \mathcal{Y}_c / N$ is a sample estimate of $\Sigma_{XY}$, the covariance between $X$ and $Y$. With this notation, we have

$$\hat{b} = \hat{\Sigma}_{XX}^{-1} \hat{\Sigma}_{XY},$$

which, by the law of large numbers, converges to $b^* = \Sigma_{XX}^{-1} \Sigma_{XY}$.

Let $a_0^* = m_Y - m_X^T b^*$. Then $f^*(x) = a_0^* + (b^*)^T x$ is the least-square optimal approximation of $Y$ by a *linear* function of $X$, and the linear predictor $\hat{f}(x) = \hat{a}_0 + \hat{b}^T x$ converges a.s. to $f^*(x)$. Of course, $f^*$ generally differs from $f : x \mapsto E(Y \mid X = x)$, which is the least-square optimal approximation of $Y$ by *any* (square-integrable) function of $X$, so that the linear estimator will have a residual bias.

### 7.1.3 Gauss-Markov theorem

If one makes the (unlikely) assumption that the linear model is exact, i.e., $f(x) = f^*(x)$, one has:

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\mathbb{E}(\hat{\beta} \mid \mathcal{X})) = \mathbb{E}((\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathbb{E}(\mathcal{Y} \mid \mathcal{X})) = \mathbb{E}((\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathcal{X} \beta) = \beta$$

and the estimator is "unbiased." Under this parametric assumption, many other properties of linear estimators can be proved, among which the well-known Gauss-Markov theorem on the optimality of least-square estimation that we now state and prove. For this theorem, for which we take (for simplicity) $q = 1$, we also assume that $\mathrm{var}(Y \mid X = x)$, the variance of $Y$ for its conditional distribution given $X$ does not depend on $x$, and denote it by $\sigma^2$. This typically correspond to the standard regression model in which one assumes that $Y = f(X) + \epsilon$ where $\epsilon$ is independent of $X$ with variance $\sigma^2$.

Recall that a symmetric matrix $A$ is said to be larger than or equal to another symmetric matrix, $B$, writing $A \geq B$, if and only if $A - B$ is positive semi-definite.

**Theorem 7.3 (Gauss-Markov)** *Assume that an estimator $\tilde{\beta}$ takes the form $\tilde{\beta} = A(\mathcal{X})\mathcal{Y}$ (it is linear) and is unbiased conditionally to $\mathcal{X}$: $\mathbb{E}_\beta(\tilde{\beta} \mid \mathcal{X}) = \beta$ (for all $\beta$). Then (under the assumptions above) the covariance matrix of $\tilde{\beta}$ cannot be smaller than that of the least square estimate, $\hat{\beta}$.*

PROOF We write $A = A(\mathcal{X})$ for short. The condition that $\mathbb{E}(A\mathcal{Y} \mid \mathcal{X}) = \beta$ for all $\beta$ yields $A\mathcal{X}\beta = \beta$ for all $\beta$, or $A\mathcal{X} = \mathrm{Id}_{\mathbb{R}^{d+1}}$ ($A$ is a $(d+1) \times N$ matrix). Since $\tilde{\beta}$ is unbiased, its covariance matrix is

$$\mathbb{E}(A\mathcal{Y}\mathcal{Y}^T A^T) - \beta\beta^T$$

and

$$\mathbb{E}(A\mathcal{Y}\mathcal{Y}^T A^T) = \mathbb{E}(\mathbb{E}(A\mathcal{Y}\mathcal{Y}^T A^T \mid \mathcal{X})) = \sigma^2 E(AA^T).$$

For $\tilde{\beta} = \hat{\beta}$, for which $A = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T$, we get $\mathbb{E}(A\mathcal{Y}\mathcal{Y}^T A^T) = \sigma^2 \mathbb{E}((\mathcal{X}^T \mathcal{X})^{-1})$. We therefore need to show that $\mathbb{E}(AA^T) \geq \mathbb{E}(\mathcal{X}^T \mathcal{X})$, i.e., that for any $u \in \mathbb{R}^{d+1}$,

$$u^T \mathbb{E}(AA^T)u \geq u^T \mathbb{E}((\mathcal{X}^T \mathcal{X})^{-1})u$$

as soon as $A\mathcal{X} = \mathrm{Id}_{\mathbb{R}^{d+1}}$. We in fact have the stronger result (without expectations):

$$A\mathcal{X} = \mathrm{Id}_{\mathbb{R}^{d+1}} \Rightarrow AA^T \geq (\mathcal{X}^T \mathcal{X})^{-1}.$$

To see this, fix $u$ and consider the problem of minimizing $F_u(A) = A \mapsto u^T AA^T u$ subject to the linear constraint $A\mathcal{X} = \mathrm{Id}_{\mathbb{R}^{d+1}}$. The Lagrange multipliers for this affine constraint can be organized in a matrix $C$ and the Lagrangian is

$$u^T AA^T u + \mathrm{trace}(C^T(A\mathcal{X} - \mathrm{Id}_{\mathbb{R}^{d+1}})).$$

Taking the derivative in $A$, we find that optimal solutions must satisfy

$$2u^T A H^T u + \text{trace}(C^T H \mathcal{X}) = 0$$

for all $H$, which yields $\text{trace}(H^T (2uu^T A + C\mathcal{X}^T)) = 0$ for all $H$. This is only possible when $2uu^T A + C\mathcal{X}^T = 0$, which in turn implies that $2uu^T A\mathcal{X} = -C\mathcal{X}^T \mathcal{X}$. Using the constraint, we get

$$C = -2uu^T (\mathcal{X}^T \mathcal{X})^{-1}$$

so that $uu^T A = uu^T (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T$. This implies that $A = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T$ (the least-square estimator) is a minimizer of $F_u(A)$ for all $u$.

Any other solution that satisfies $uu^T A = uu^T (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T$ for all $u$. Taking $u = \mathfrak{e}_i$ and summing over $i$ (with $\sum_{i=1}^{d+1} \mathfrak{e}_i \mathfrak{e}_i^T = \text{Id}_{\mathbb{R}^{d+1}}$) yields $A = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T$.  ∎

### 7.1.4   Kernel Version

We now assume that $X$ takes its values in an arbitrary set $\mathcal{R}_X$, with a representation $h : \mathcal{R}_X \to H$ into an inner-product space. This representation does not need to be explicit or computable, but the associated kernel $K(x, y) = \langle h(x), h(y) \rangle_H$ is assumed to be known and easy to compute. (Recall that, from chapter 6, a positive kernel is always associated with an inner-product space.) In particular, any algorithm in this context should only rely on the kernel, and the function $h$ only has a conceptual role.

Assume that $q = 1$ to lighten the notation, so that the dependent variable is scalar-valued. We here let the space of predictors be

$$\mathcal{F} = \{f(x) = a_0 + \langle b, h(x) \rangle_H : a_0 \in \mathbb{R}, b \in H\}.$$

The residual sum of squares associated with this function space is

$$RSS(a_0, b) = \sum_{k=1}^{N} (y_k - a_0 - \langle b, h(x_k) \rangle)^2.$$

The following result (or results similar to it) is a key step in almost all kernel methods in machine learning.

**Proposition 7.4** *Let $V = \text{span}(h(x_1), \ldots, h(x_N))$ be the finite-dimensional subspace of $H$ generated by the feature functions evaluated on training input data. Then*

$$RSS(a_0, b) = RSS(a_0, \pi_V(b)).$$

*where $\pi_V$ is the orthogonal projection on $V$.*

PROOF The justification is immediate: since $h(x_k) \in V$, we have

$$\langle b, h(x_k) \rangle_H = \langle \pi_V(b), h(x_k) \rangle_H$$

for all $b \in H$. ∎

This shows that there is no loss of generality in restricting the minimization of the residual sum of squares to $b \in V$. Such a $b$ takes the form

$$b = \sum_{k=1}^{N} \alpha_k h(x_k) \qquad (7.3)$$

and the regression problem can be reformulated as a function of the coefficients $\alpha_1, \ldots, \alpha_N \in \mathbb{R}$, with

$$f(x) = a_0 + \sum_{k=1}^{N} \alpha_k \langle h(x), h(x_k) \rangle_H = a_0 + \sum_{k=1}^{N} \alpha_k K(x, x_k),$$

which only depends on the kernel. (This reduction is often referred to as the "kernel trick.")

However, the solution of the problem is, in this context, not very interesting. Indeed, assume that $K$ is positive definite and that all observations in the training set are distinct. Then the matrix $\mathcal{K}(x_1, \ldots, x_N)$ formed by the kernel evaluations $K(x_i, x_j)$ is invertible, and one can solve exactly the equations

$$y_k = \sum_{j=1}^{N} \alpha_j K(x_k, x_j), \quad k = 1, \ldots, N$$

to get a zero RSS with $a_0 = 0$. Unless there is no noise, such a solution will certainly overfit the data. If $K$ is not positive definite, and the dimension of $V$ is less than $N$ (since this would place us in the previous situation otherwise), then it is more efficient to work directly in a basis of $V$ rather than using the over-parametrized kernel representation. We will see however, starting with the next section, that kernel methods become highly relevant as soon as the regression is estimated with some control on the size of the regression coefficients, $b$.

## 7.2 Ridge regression and Lasso

### 7.2.1 Ridge Regression

**Method.** When the set $\mathcal{F}$ of possible predictors is too large, some additional complexity control is needed to reduce the estimation variance. One simple approach

is to limit the number of parameters to be estimated, which, for regression, corresponds to limiting the number of possible predictors. This is related to the methods of Sieves mentioned in section 4.1. In contrast, ridge regression and lasso control the size of the parameters, as captured by their norm.

In both cases, one assigns a measure of complexity, denoted $f \mapsto \gamma(f) \geq 0$, to each element $f \in \mathcal{F}$. Given $\gamma$, one can either optimize this predictor (using, for example, the RSS) with the constraint that $\gamma(f) \leq C$ for some constant $C$, or add a penalty $\lambda \gamma(f)$ to the objective function for some $\lambda > 0$. In general, the two approaches (constraint or penalty) are equivalent.

In linear spaces, complexity measures are often associated with a norm, and ridge regression uses the sum of squares of coefficients of the prediction matrix $b$, minimizing

$$\sum_{k=1}^{N} |y_k - a_0 - b^T x_k|^2 + \lambda \operatorname{trace}(b^T b), \tag{7.4}$$

which can be written in vector form as

$$|\mathcal{Y} - \mathcal{X}\beta|_2^2 + \lambda \operatorname{trace}(\beta^T \Delta \beta),$$

where $\Delta = \operatorname{diag}(0, 1, \ldots, 1)$. In the following, we will work with an unspecified $(d + 1) \times (d + 1)$ symmetric positive semi-definite matrix $\Delta$. Various choices are indeed possible, for example, $\Delta = \operatorname{diag}(0, \hat{\sigma}^2(1), \ldots, \hat{\sigma}^2(d))$, where $\hat{\sigma}^2(i)$ is the empirical variance of the $i$th coordinate of $X$ in the training set. This last choice is quite natural, because it ensures that, whenever one of the variable $X^{(i)}$ is rescaled by a factor $c$, the corresponding optimal $i^{\text{th}}$ row of $b^T$ is rescaled by $1/c$, leaving the predictor unchanged.

Under this assumption, the optimal parameter is

$$\hat{\beta}^{\lambda} = (\mathcal{X}^T \mathcal{X} + \lambda \Delta)^{-1} \mathcal{X}^T \mathcal{Y},$$

with a proof similar to that made for least-square regression. We obviously retrieve the original formula for regression when $\lambda = 0$.

Alternatively, assuming that $\Delta = \begin{pmatrix} 0 & 0 \\ 0 & \Delta' \end{pmatrix}$, so that no penalty is imposed on the intercept, we have

$$\hat{b}^{\lambda} = (\mathcal{X}_c^T \mathcal{X}_c + \lambda \Delta')^{-1} \mathcal{X}_c^T \mathcal{Y}_c \tag{7.5}$$

and $\hat{a_0}^{\lambda} = \bar{y} - (\hat{b}^{\lambda})^T \bar{x}$. The proof of these statements is left to the reader.

**Analysis in a special case** To illustrate the impact of the penalty term on balancing bias and variance, we now make a computation in the special case when $Y = \tilde{X}\beta + \epsilon$, where $\text{var}(\epsilon) = \sigma^2$ and $\epsilon$ is independent of $X$. In the following computation, we assume that the training set is fixed (or rather, compute probabilities and expectations conditionally to it). Also, to simplify notation, we denote

$$S_\lambda = \mathcal{X}^T\mathcal{X} + \lambda\Delta = \sum_{k=1}^{N} \tilde{x}_k^T \tilde{x}_k + \lambda\Delta$$

and $\Sigma = E(\tilde{X}^T\tilde{X})$ for a single realization of $X$. Finally, we assume that $q = 1$, also to simplify the discussion.

The mean-square prediction error is

$$
\begin{aligned}
R(\lambda) &= E((Y - \tilde{X}^T\hat{\beta}_\lambda)^2) \\
&= E((\tilde{X}^T(\beta - \hat{\beta}_\lambda) + \epsilon)^2) \\
&= (\hat{\beta}_\lambda - \beta)^T\Sigma(\hat{\beta}_\lambda - \beta) + \sigma^2.
\end{aligned}
$$

Denote by $\epsilon_k$ the (true) residual $\epsilon_k = y_k - \tilde{x}_k^T\beta$ on training data and by $\epsilon$ the vector stacking these residuals. We have, writing $S_0 = S_\lambda - \lambda\Delta$,

$$
\begin{aligned}
\hat{\beta}_\lambda &= S_\lambda^{-1}\mathcal{X}^T\mathcal{Y} \\
&= S_\lambda^{-1}S_0\beta + S_\lambda^{-1}\mathcal{X}^T\epsilon \\
&= \beta - \lambda S_\lambda^{-1}\Delta\beta + S_\lambda^{-1}\mathcal{X}^T\epsilon
\end{aligned}
$$

So we can rewrite

$$R(\lambda) = \lambda^2\beta^T\Delta S_\lambda^{-1}\Sigma S_\lambda^{-1}\Delta\beta - 2\lambda\epsilon^T\mathcal{X}S_\lambda^{-1}\Sigma S_\lambda^{-1}\Delta\beta + \epsilon^T\mathcal{X}S_\lambda^{-1}\Sigma S_\lambda^{-1}\mathcal{X}^T\epsilon + \sigma^2.$$

Let us analyze the quantities that depend on the training set in this expression. The first one is $S_\lambda = S_0 + \lambda\Delta$. From the law of large numbers, $S_0/N \to \Sigma$ when $N$ tends to infinity, so that, assuming in addition that $\lambda = \lambda_N = O(N)$, we have $S_\lambda^{-1} = O(1/N)$. The second one is

$$\epsilon^T\mathcal{X} = \sum_{k=1}^{N} \epsilon_k\tilde{x}_k$$

which, according to the central limit theorem, is such that

$$N^{-1/2}\epsilon^T\mathcal{X} \sim \mathcal{N}(0, \sigma^2\text{Var}(\tilde{X}))$$

when $N \to \infty$. So, we can expect the coefficient of $\lambda^2$ in $R(\lambda)$ to have order $N^{-2}$, the coefficient of $\lambda$ to have order $N^{-3/2}$ and the constant coefficient of have order $N^{-1}$. This suggests taking $\lambda = \mu\sqrt{N}$ so that all coefficients have roughly the same order when expanding in powers of $\mu$.

This gives $S_\lambda = N(S_0/N + \mu\Delta/\sqrt{N}) \simeq N\Sigma$ and we make the approximation, letting $\xi = N^{-1/2}\sigma^{-1/2}\epsilon\mathcal{X}^T$ and $\gamma = \Sigma^{-1/2}\Delta\beta$, that

$$N(R(\lambda) - \sigma^2) \simeq \mu^2|\gamma|^2 - 2\mu\xi^T\gamma + \xi^T\xi.$$

With this approximation, the optimal $\mu$ should be

$$\mu = \frac{\xi^T\gamma}{|\gamma|^2}.$$

Of course, this $\mu$ cannot be computed from data, but we can see that, since $\xi$ converges to a centered Gaussian random variable, its value cannot be too large. It is therefore natural to choose $\mu$ to be constant and use ridge regression in the form

$$\sum_{k=1}^{N}(y_k - \tilde{x}_k^T\beta)^2 + \sqrt{N}\mu\beta^T\Delta\beta.$$

In all cases, the mere fact that we find that the optimal $\mu$ is not 0 shows that, under the simplifying (and optimistic) assumptions that we made for this computation, allowing for a penalty term always reduces the prediction error. In other terms, introducing some estimation bias in order to reduce the variance is beneficial.

**Kernel Ridge Regression**   We now return to the feature-space situation and take $h : \mathcal{R}_X \to H$ with associated kernel $K$. We still take $q = 1$ for simplicity. One formulates the ridge regression problem in this context as the minimization of

$$\sum_{k=1}^{N}(y_l - a_0 - \langle b, h(x_l)\rangle_H)^2 + \lambda\|b\|_H^2$$

with respect to $\beta = (a_0, b)$. Introducing the space $V$ generated by the feature function evaluated on the training set, we know from proposition 7.4 that replacing $b$ by $\pi_V(b)$ leaves the residual sum of squares invariant. Moreover, one has $\|\pi_V(b)\|_H^2 \leq \|b\|_H^2$ with equality if and only if $b \in V$. This shows that the solution $b$ must belong to $V$ and therefore take the form (7.3).

Using this expression, one finds that the problem is reduced to finding the minimum of

$$\sum_{k=1}^{N}\left(y_k - a_0 - \sum_{l=1}^{N}K(x_l,x_k)\alpha_l\right)^2 + \lambda\sum_{k,l=1}^{N}\alpha_k\alpha_l K(x_k,x_l)$$

with respect to $a_0, \alpha_1, \ldots, \alpha_N$. Recall that we have denoted by $\mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$ the kernel matrix with entries $K(x_i, x_j)$, $i, j = 1, \ldots, N$. We will assume in the following that $\mathcal{K}$ is invertible.

Introduce the vector $\mathbb{1}_N \in \mathbb{R}^N$ with all coordinates equal to one. Let

$$\tilde{\mathcal{K}} = \begin{pmatrix} \mathbb{1}_N & \mathcal{K} \end{pmatrix} \text{ and } \mathcal{K}' = \begin{pmatrix} 0 & 0 \\ 0 & \mathcal{K} \end{pmatrix}.$$

Let $\alpha \in \mathbb{R}^N$ be the vector with coefficients $\alpha_1, \ldots, \alpha_N$ and $\tilde{\alpha} = \begin{pmatrix} a_0 \\ \alpha \end{pmatrix}$. With this notation, the function to minimize is

$$F(\alpha) = |\mathcal{Y} - \tilde{\mathcal{K}}\tilde{\alpha}|^2 + \lambda \tilde{\alpha}^T \mathcal{K}' \tilde{\alpha}.$$

This takes the same form as standard ridge regression, replacing $\beta$ by $\tilde{\alpha}$, $\mathcal{X}$ by $\tilde{\mathcal{K}}$ and $\Delta$ by $\mathcal{K}'$. The solution therefore is

$$\tilde{\alpha}^\lambda = (\tilde{\mathcal{K}}^T \tilde{\mathcal{K}} + \lambda \mathcal{K}')^{-1} \tilde{\mathcal{K}}^T \mathcal{Y}.$$

Note that $\mathcal{K}$ being invertible implies that $\tilde{\mathcal{K}}^T \tilde{\mathcal{K}} + \lambda \mathcal{K}'$ is invertible. [1]

To write the equivalent of (7.5), we need to use the equivalent of the matrix $\mathcal{X}_c$, that is, the matrix $\mathcal{K}$ with the average of the $j$th column subtracted to each $(i, j)$ entry, given by:

$$\mathcal{K}_c = \mathcal{K} - \frac{1}{N} \mathbb{1}_N \mathbb{1}_N^T \mathcal{K}.$$

Introduce the matrix $P = \text{Id} - \mathbb{1}_N \mathbb{1}_N^T / N$. It is easily checked that $P^2 = P$ ($P$ is a projection matrix). Since $\mathcal{K}_c = P\mathcal{K}$, we have $\mathcal{K}_c^T \mathcal{K}_c = \mathcal{K}P\mathcal{K}$. One deduces from this the expression of the optimal vector $\alpha^\lambda$, namely,

$$\alpha^\lambda = (\mathcal{K}P\mathcal{K} + \lambda \mathcal{K})^{-1} \mathcal{K}P\mathcal{Y}_c = (P\mathcal{K} + \lambda \text{Id}_{\mathbb{R}^N})^{-1} \mathcal{Y}_c$$

where we have, in addition, used the fact that $P\mathcal{Y}_c = \mathcal{Y}_c$. Finally, the intercept is given by

$$a_0 = \bar{y} - \frac{1}{N} (\alpha^\lambda)^T K \mathbb{1}_N.$$

### 7.2.2 Equivalence of constrained and penalized formulations

**Case of ridge regression.** Returning to the basic case (without feature space), we now introduce an alternate formulation of ridge regression. Let *ridge*($\lambda$) denote the ridge regression problem that we have considered so far, for some parameter

---

[1] Indeed, let $u = \begin{pmatrix} w_0 \\ w \end{pmatrix}$ with $w_0 \in \mathbb{R}$ and $w \in \mathbb{R}^N$ be such that $u^T(\tilde{\mathcal{K}}^T \tilde{\mathcal{K}} + \lambda \mathcal{K}')u = 0$. This requires $\tilde{\mathcal{K}}u = 0$ and $u^T \mathcal{K}'u = 0$. The latter quantity is $w^T \mathcal{K}w$, which shows that $w = 0$ since $\mathcal{K}$ has rank $N$. Then $\tilde{\mathcal{K}} = \mathbb{1}_N w_0$ so that $w_0 = 0$ also.

$\lambda$. Consider now the following problem, which will be called $ridge'(C)$: minimize $\sum_{k=1}^{N} |y_k - \tilde{x}_k^T \beta|^2$ subject to the constraint $\beta^T \Delta \beta \leq C$. We claim that this problem is equivalent to the ridge regression problem, in the following sense: for any $C$, there exists a $\lambda$ such that the solution of $ridge'(C)$ coincides with the solution of $ridge(\lambda)$ and vice-versa.

Indeed, fix a $C > 0$. Consider an optimal $\beta$ for $ridge'(C)$. Assuming as above that $\Delta$ is symmetric positive semi-definite, we let $V$ be its null space and $P_V$ the orthogonal projection on $V$. Write $\beta = \beta_1 + \beta_2$ with $\beta_1 = P_V \beta$. Let $d_1$ and $d_2$ be the respective dimensions of $V$ and $V^\perp$ so that $d_1 + d_2 = d$. Identifying $\mathbb{R}^d$ with the product space $V \times V^\perp$ (i.e., making a linear change of coordinates), the problem can be rewritten as the minimization of

$$|\mathcal{Y} - \mathcal{X}_1 \beta_1 - \mathcal{X}_2 \beta_2|^2$$

subject to $\beta_2^T \Delta \beta_2 \leq C$, where $\mathcal{X}_1$ (resp. $\mathcal{X}_2$) is $N \times d_1$ (resp. $N \times d_2$).

The gradient of the constraint $\gamma(\beta_2) = \beta_2^T \Delta \beta_2 - C$ is $\nabla \gamma(\beta_2) = 2\Delta \beta_2$. Assume first that $\Delta \beta_2 \neq 0$. Then the solution must satisfy the KKT conditions, which require that there exists $\mu \geq 0$ such that $\beta$ is a stationary point of the Lagrangian

$$|\mathcal{Y} - \mathcal{X}_1 \beta_1 - \mathcal{X}_2 \beta_2|^2 + \mu \beta_2^T \Delta \beta_2,$$

with $\mu > 0$ only possible if $\beta^T \Delta \beta = C$. This requires that

$$\mathcal{X}_1^T \mathcal{X}_1 \beta_1 + \mathcal{X}_1^T \mathcal{X}_2 \beta_2 = \mathcal{X}^T \mathcal{Y},$$
$$\mathcal{X}_2^T \mathcal{X}_1 \beta_1 + \mathcal{X}_2^T \mathcal{X}_2 \beta_2 + \mu \Delta \beta_2 = \mathcal{X}^T \mathcal{Y}.$$

Since $\Delta \beta_1 = 0$, and using $\mathcal{X} = (\mathcal{X}_1, \mathcal{X}_2)$, we have

$$\beta = (\mathcal{X}^T \mathcal{X} + \mu \Delta)^{-1} \mathcal{X}^T Y,$$

which is the only solution of ridge$(\mu)$.

If $\Delta \beta_2 = 0$, then, necessarily, $\beta_2 = 0$. Since $C > 0$, $\beta$ must then be the solution of the unconstrained problem, which is ridge$(0)$.

Conversely, any solution $\beta$ of $ridge(\lambda)$ satisfies the first-order optimality conditions for $ridge'(C)$ for $C = \beta^T \Delta \beta$ (or any $C \geq \beta^T \Delta \beta$ if $\lambda = 0$). This shows the equivalence of the two problems.

**General case.**    We now consider this equivalence in a more general setting. Consider a penalized optimization problem, denoted $var(\lambda)$ which consists in minimizing in $\beta$ some objective function of the form $U(\beta) + \lambda \varphi(\beta), \lambda \geq 0$. Consider also the family of problems $var'(C)$, with $C > \inf(\varphi)$, which minimize $U(\beta)$ subject to $\varphi(\beta) \leq C$.

We make the following assumptions.

(i) $U$ and $\varphi$ are continuous functions from $\mathbb{R}^n$ to $\mathbb{R}$.

(ii) $\varphi(\beta) \to \infty$ when $\beta \to \infty$.

(iii) For any $\lambda \geq 0$, there is a unique solution of $var(\lambda)$, denoted $\beta_\lambda$.

(iv) For any $C$, there is a unique solution of $var'(C)$. denoted $\beta'_C$.

Assumptions (ii) and (iv) are true, in particular, when $U$ is strictly convex, $\varphi$ is convex and $U$ has compact level sets. We show that, with these assumptions, the two families of problems are equivalent.

We first discuss the penalized problems and prove the following proposition, which has its own interest.

**Proposition 7.5** *The function* $\lambda \mapsto U(\beta_\lambda)$ *is nondecreasing, and* $\lambda \mapsto \varphi(\beta_\lambda)$ *is non-increasing, with*

$$\lim_{\lambda \to \infty} \varphi(\beta_\lambda) = \inf(\varphi).$$

*Moreover,* $\beta_\lambda$ *varies continuously as a function of* $\lambda$.

PROOF Consider two parameters $\lambda$ and $\lambda'$. We have

$$U(\beta_\lambda) + \lambda\varphi(\beta_\lambda) \leq U(\beta_{\lambda'}) + \lambda\varphi(\beta_{\lambda'})$$
$$\text{and } U(\beta_{\lambda'}) + \lambda'\varphi(\beta_{\lambda'}) \leq U(\beta_\lambda) + \lambda'\varphi(\beta_\lambda)$$

since both left-hand sides are minimizers. This implies

$$\lambda(\varphi(\beta_\lambda) - \varphi(\beta_{\lambda'})) \leq U(\beta_{\lambda'}) - U(\beta_\lambda) \leq \lambda'(\varphi(\beta_\lambda) - \varphi(\beta_{\lambda'})). \tag{7.6}$$

In particular: $(\lambda' - \lambda)(\varphi(\beta_\lambda) - \varphi(\beta_{\lambda'})) \geq 0$. Assume that $\lambda < \lambda'$. Then this last inequality implies $\varphi(\beta_\lambda) \geq \varphi(\beta_{\lambda'})$ and (7.6) then implies that $U(\beta_\lambda) \leq U(\beta_{\lambda'})$, which proves the first part of the proposition.

Now assume that there exists $\epsilon > 0$ such that $\varphi(\beta_\lambda) > \inf \varphi + \epsilon$ for all $\lambda \geq 0$. Take $\tilde{\beta}$ such that $\varphi(\tilde{\beta}) \leq \inf \varphi + \epsilon/2$. For any $\lambda > 0$, we have

$$U(\beta_\lambda) + \lambda\varphi(\beta_\lambda) \leq U(\tilde{\beta}) + \lambda\varphi(\tilde{\beta})$$

so that $U(\beta_\lambda) < U(\tilde{\beta}) - \lambda\epsilon/2$. Since $U(\beta_\lambda) \geq U(a_0)$, we get $U(a_0) = -\infty$, which is a contradiction. This shows that $\varphi(\beta_\lambda)$ tends to $\inf(\varphi)$ when $\lambda$ tends to infinity.

We now prove that $\lambda \mapsto \beta_\lambda$ is continuous. Define $G(\beta, \lambda) = U(\beta) + \lambda\varphi(\beta)$. Since we assume that $\varphi(\beta) \to \infty$ when $\beta \to \infty$, and we have just proved that $\varphi(\beta_\lambda) \leq \varphi(a_0)$ for any $\lambda$, we obtain the fact that the set $(|\beta_\lambda|, \lambda \geq 0)$ is bounded, say by a constant $B \geq 0$.

Consider a sequence $\lambda_n$ that converges to $\lambda$. We want to prove that $\beta_{\lambda_n} \to \beta_\lambda$, for which (because $\beta_\lambda$ is bounded) it suffices to show that if any subsequence of $(\beta_{\lambda_n})$ converges to some $\tilde{\beta}$, then $\tilde{\beta} = \beta_\lambda$.

So, consider such a converging subsequence, that we will still denote by $\beta_{\lambda_n}$ for convenience. Since $G$ is continuous, one has $G(\beta_{\lambda_n}, \lambda_n) \to G(\tilde{\beta}, \lambda)$ when $n$ tends to infinity. Let us prove that $G(\beta_\lambda, \lambda)$ is continuous in $\lambda$. For any pair $\lambda, \lambda'$ and any $\beta$, we have

$$G(\beta_{\lambda'}, \lambda') \leq G(\beta_\lambda, \lambda') = G(\beta_\lambda, \lambda) + (\lambda' - \lambda)\varphi(\beta_\lambda) \leq G(\beta_\lambda, \lambda) + |\lambda' - \lambda|\varphi(a_0).$$

This yields, by symmetry, $|G(\beta_{\lambda'}, \lambda') - G(\beta_\lambda, \lambda)| \leq \varphi(a_0)|\lambda - \lambda'|$, proving the continuity in $\lambda$.

So we must have $G(\tilde{\beta}, \lambda) = G(\beta_\lambda, \lambda)$. This implies that both $\tilde{\beta}$ and $\beta_\lambda$ are solutions of $var(\lambda)$, so that $\beta_\lambda = \tilde{\beta}$ because we assume that the solution is unique.    ∎

We now prove that the classes of problems $var(\lambda)$ and $var'(C)$ are equivalent. First, $\beta_\lambda$ is a minimizer of $U(\beta)$ subject to the constraint $\varphi(\beta) \leq C$, with $C = \varphi(\beta_\lambda)$. Indeed, if $U(\beta) < U(\beta_\lambda)$ for some $\beta$ with $\varphi(\beta) \leq \varphi(\beta_\lambda)$, then $U(\beta) + \lambda\varphi(\beta) < U(\beta_\lambda) + \lambda\varphi(\beta_\lambda)$ which is a contradiction. So $\beta_\lambda = \beta'_{\varphi(\beta_\lambda)}$. Using the continuity of $\beta_\lambda$ and $\varphi$, this proves the equivalence of the problems when $C$ is in the interval $(a, \varphi(a_0))$ where $a = \lim_{\lambda\to\infty} \varphi(\beta_\lambda) = \inf(\varphi)$.

So, it remains to consider the case $C > \varphi(a_0)$. For such a $C$, the solution of $var'(C)$ must be $a_0$ since it is a solution of the unconstrained problem, and satisfies the constraint.

### 7.2.3   Lasso regression

**Problem statement**   Assume that the output variable is scalar, i.e., $q = 1$. Let $\hat{\sigma}^2(i)$ be the empirical variance of the $i$th variable $X^{(i)}$. Then, the lasso estimator is defined as a minimizer of $\sum_{k=1}^N (y_k - \tilde{x}_k^T \beta)^2$ subject to the constraint $\sum_{i=1}^d \hat{\sigma}(i)|\beta^{(i)}| \leq C$. Compared to ridge regression, the sum of squares for $\beta$ is simply replaced by a weighted sum of absolute values, but we will see that this change may significantly affect the nature of the solutions.

As we have just seen, the penalized formulation, minimizing

$$\sum_{k=1}^N (y_k - \tilde{x}_k^T \beta)^2 + \lambda \sum_{i=1}^d \hat{\sigma}(i)|\beta^{(i)}|$$

provides an equivalent family of problems, on which we will focus (because it is easier to analyze). Since one uses a non-Euclidean norm in the penalty, there is no

kernel version of the lasso and we only discuss the method in the original input space $\mathcal{R} = \mathbb{R}^d$.

For a vector $a \in \mathbb{R}^k$, we let $|a|_1 = |a^{(1)}| + \cdots + |a^{(k)}|$, the $\ell^1$ norm of $a$. Using the previous notation for $\mathcal{Y}$ and $\mathcal{X}$, the quantity to minimize can be rewritten as

$$|\mathcal{Y} - \mathcal{X}\beta|^2 + \lambda |D\beta|_1$$

where $D$ is the $d \times (d+1)$ matrix with $d(i, i+1) = \hat{\sigma}(i)$ for $i = 1, \ldots, d$ and all other coefficients equal to 0. This is a convex optimization problem which, unlike ridge regression, does not have a closed form solution.

**ADMM.** The alternating direction method of multipliers (ADMM) that was described in section 3.6, (3.59) is one of the state-of-the-art algorithm to solve the lasso problem, especially in large dimensions. Other iterative methods include subgradient descent (see the example in section 3.5.4) and proximal gradient descent. Since $x$ has a different meaning here, we change the notation in (3.59) by replacing $x, z, u$ by $\beta, \gamma, \tau$, and rewrite the lasso problem as the minimization of

$$|\mathcal{Y} - \mathcal{X}\beta|^2 + \lambda |\gamma|_1$$

subject to $D\beta - \gamma = 0$. Applying (3.59) with $A = D$, $B = -\text{Id}$ and $c = 0$, the ADMM iterations are

$$\begin{cases} \beta(n+1) & = \text{argmin}_\beta \left( |\mathcal{Y} - \mathcal{X}\beta|^2 + \dfrac{1}{2\rho} |D\beta - \gamma(n) + \tau(n)|^2 \right) \\[2ex] \gamma^{(i)}(n+1) & = \text{argmin}_t \left( \lambda |t| + \dfrac{1}{2\rho}(t - D\beta^{(i)}(n+1) - \tau^{(i)}(n))^2 \right), \ i = 1, \ldots, d \\[2ex] \tau(n+1) & = \tau(n) + D\beta(n+1) - \gamma(n+1) \end{cases}$$

The solutions of both minimization problems are explicit, yielding the following algorithm, which converges to a solution if $\rho$ is small enough.

---

**Algorithm 7.1 (ADMM for lasso)**
Let $\rho > 0$ be chosen. Starting with initial values $\beta^{(0)}$, $\gamma^{(0)}$ and $\tau^{(0)}$, the ADMM algorithm for lasso iterates:

$$\begin{cases} \beta(n+1) & = \left( \mathcal{X}^T \mathcal{X} + \dfrac{D^T D}{2\rho} \right)^{-1} \left( \mathcal{X}^T \mathcal{Y} + \dfrac{D^T}{2\rho}(\gamma(n) - \tau(n)) \right) \\[2ex] \gamma^{(i)}(n+1) & = S_{\lambda\rho} \left( D\beta^{(i)}(n+1) + \tau^{(i)}(n) \right), \ i = 1, \ldots, d \\[2ex] \tau(n+1) & = \tau(n) + D\beta(n+1) - \gamma(n+1) \end{cases}$$

until the difference between the variables at steps $n$ and $n + 1$ is below a small toler-ance level. Here, $S_{\lambda\rho}$ is the so-called *shrinkage operator*

$$S_{\lambda\rho}(v) = \begin{cases} v - \lambda\rho & \text{if } v \geq \lambda\rho \\ 0 & \text{if } |v| \leq \lambda\rho \\ v + \lambda\rho & \text{if } v \leq -\lambda\rho \end{cases}$$

Note that the ADMM algorithm makes an iterative approximation of the constraints, so that they are only satisfied at some precision level when the algorithm is stopped.

**Exact computation.**   We now provide a more detailed characterization of the solu-tion of the lasso problem and analyze, in particular, how this solution changes when $\lambda$ (or $C$) varies. To simplify the exposition, and without loss of generality, we will as-sume that the variables have been normalized so that $\hat{\sigma}(i) = 1$ and the penalty simply is the sum of absolute values. Let

$$G_\lambda(\beta) = \sum_{k=1}^{N}(y_k - a_0 - x_k^T b)^2 + \lambda \sum_{i=1}^{d}|b(i)|.$$

The following proposition, in which we let

$$r_b = \frac{1}{N}\sum_{k=1}^{N}(y_k - a_0 - x_k^T b)x_k,$$

characterizes the solution of the lasso.

**Proposition 7.6** *The pair $(a_0, b)$ is the optimal solution of the lasso problem with param-eter $\lambda$ if and only if $a_0 = \bar{y} - \bar{x}^T b$ and, for all $i = 1, \ldots, d$,*

$$|r_b^{(i)}| \leq \frac{\lambda}{2N} \tag{7.7}$$

*with*

$$r_b^{(i)} = \text{sign}(b^{(i)})\frac{\lambda}{2N} \text{ if } b^{(i)} \neq 0. \tag{7.8}$$

*In particular $|r_b^{(i)}| < \lambda/(2N)$ implies $b^{(i)} = 0$.*

PROOF Using the subdifferential calculus in theorem 3.45, one can compute the sub-gradients of $G$ by adding the subdifferentials of the terms that compose it. All these terms are differentiable except $|b^{(i)}|$ when $b^{(i)} = 0$, and the subdifferential of $t \mapsto |t|$ at $t = 0$ is the interval $[-1, 1]$.

This shows that $g \in \partial G_\lambda(\beta)$ if and only if

$$g = -2N r_b + \lambda z$$

with $z^{(i)} = \text{sign}(b^{(i)})$ if $b^{(i)} \neq 0$ and $|z^{(i)}| \leq 1$ otherwise. Proposition 7.6 immediately follows by taking $g = 0$. $\blacksquare$

Let $\zeta = \text{sign}(b)$, the vector formed by the signs of the coordinates of $b$, with $\text{sign}(0) = 0$. Then proposition 7.6 uniquely specifies $a_0$ and $b$ once $\lambda$ and $\zeta$ are known. Indeed, let $J = J_\zeta$ denote the ordered subset of indices $j \in \{1, \ldots, d\}$ such that $\zeta^{(j)} \neq 0$, and let $b(J)$, $x_k(J)$, $\zeta(J)$, etc., denote the restrictions of vectors to these indices. Equation (7.8) can be rewritten as (after replacing $a_0$ by its optimal value)

$$\mathcal{X}_c(J)^T \mathcal{X}_c(J) b(J) = \mathcal{X}_c(J)^T \mathcal{Y}_c - \frac{\lambda}{2} \zeta(J)$$

where

$$\mathcal{X}_c(J) = \begin{pmatrix} (x_1(J) - \overline{x}(J))^T \\ \vdots \\ (x_N(J) - \overline{x}(J))^T \end{pmatrix}.$$

This yields

$$b(J) = (\mathcal{X}_c(J)^T \mathcal{X}_c(J))^{-1} \left( \mathcal{X}_c(J)^T \mathcal{Y}_c - \frac{\lambda}{2} \zeta(J) \right), \tag{7.9}$$

which fully determine $b$ since $b^{(j)} = 0$ if $j \notin J$, by definition.

For given $\lambda$, only one sign configuration $\zeta$ will provide the correct solution, with correct signs for nonzero values of $b$ above, and correct inequalities on $r_b$. Calling this configuration $\zeta_\lambda$, one can note that if $\zeta_\lambda$ is known for a given value of $\lambda$, it remains valid if we increase or decrease $\lambda$ until one of the optimality conditions changes, i.e., either one of the coordinates $b^{(i)}, i \in J_{\zeta_\lambda}$, vanishes, or one of the inequalities for $i \notin J_{\zeta_\lambda}$ becomes an equality. Moreover, proposition 7.6 shows that between these events both $b$ and therefore $r_b$ depend linearly on $\lambda$, which makes easy the task of determining maximal intervals around a given $\lambda$ over which $\zeta$ remains unchanged.

Note that solutions are known for $\lambda = 0$ (standard least squares) and for $\lambda$ large enough (for which $b = 0$). Indeed, for $b = 0$ to be a solution, it suffices that

$$\lambda > \lambda_0 \overset{\Delta}{=} 2 \max_i \left| \sum_{k=1}^{N} (y_k - \overline{y})(x_k^{(i)} - \overline{x}^{(i)}) \right|.$$

These remarks set the stage for an algorithm computing the optimal solution of the lasso problem for all values of $\lambda$, starting either from $\lambda = 0$ or $\lambda > \lambda_0$. We will describe this algorithm starting for $\lambda > \lambda_0$, which has the merit to avoid complications

due to underconstrained least squares when $d$ is large. For this purpose, we need a little more notation. For a given $\zeta$, let

$$b_\zeta = (\mathcal{X}_c(J_\zeta)^T \mathcal{X}_c(J_\zeta))^{-1} \mathcal{X}_c(J_\zeta)^T \mathcal{Y}_c$$

and

$$u_\zeta = \frac{1}{2}(\mathcal{X}_c(J_\zeta)^T \mathcal{X}_c(J_\zeta))^{-1} \zeta(J_\zeta),$$

so that $b(J_\zeta) = b_\zeta - \lambda u_\zeta$. The residuals then take the form

$$r_b^{(i)} = \frac{1}{N} \sum_{k=1}^{N} (y_k - a_0 - b_\zeta^T x_k) x_k^{(i)} + \frac{\lambda}{N} \sum_{k=1}^{N} (x_k^T u_\zeta)(x_k^{(i)} - \overline{x}^{(i)})$$

$$= \rho_\zeta^{(i)} + \lambda d_\zeta^{(i)},$$

where the last equation defines $\rho_\zeta$ and $d_\zeta$.

Assume that one wants to minimize $G_{\lambda^*}$ for some $\lambda^* > 0$. We need to describe the sequence of changes to the minimizers of $G_\lambda$ when $\lambda$ decreases from some value larger than $\lambda_0$ to the value $\lambda^*$.

If $\lambda^* \geq \lambda_0$, then the optimal solution is $b = 0$, so we can assume that $\lambda^* < \lambda_0$. When $\lambda$ is slightly smaller than $\lambda_0$, one needs to introduce some non-zero values in $\zeta$. Those values are at the indexes $i$ such that

$$\lambda_0 = 2 \left| \sum_{k=1}^{N} (y_k - \overline{y})(x_k^{(i)} - \overline{x}^{(i)}) \right|$$

The sign of $\zeta^{(i)}$ is also determined since $\text{sign}(b^{(i)}) = \text{sign}(r_b^{(i)})$ when $b^{(i)} \neq 0$.

The algorithm will then continue by progressively adding non-zero entries to $\zeta$ when the covariance between some unused variables and the residual becomes too large, or by removing non-zero values when the optimal $b$ crosses a zero. We now describe it in detail.

---

**Algorithm 7.2 (Exact minimization for lasso)**

1. Initialization: let $\lambda(0) = 1 + \lambda_0$, $\sigma(0) = 0$ and the corresponding values $a_0(0) = \overline{y}$ and $b(0) = 0$.

2. Assume that the algorithm has reached step $n$ with current variables $\lambda(n)$, $\sigma(n)$, $a_0(n)$ and $b(n)$.

3. Determine the first $\lambda' < \lambda(n)$ for which either

   (i) For some $i$, $\zeta^{(i)}(n) \neq 0$ and $b_{\zeta(n)}^{(i)} - \lambda' u_{\zeta(n)}^{(i)} = 0$.

(ii) For some $i$, $\zeta^{(i)}(n) = 0$ and $(1 - 2Nd_{\zeta(n)}^{(i)})\lambda' - 2N\rho_{\zeta(n)} = 0$.

(iii) For some $i$, $\zeta^{(i)}(n) = 0$ and $(1 + 2Nd_{\zeta(n)}^{(i)})\lambda' + 2N\rho_{\zeta(n)} = 0$.

4. Then, there are two cases:

(a) If $\lambda' \geq \lambda^*$, set $\lambda(n+1) = \lambda'$. Let $\zeta^{(i)}(n+1) = \zeta^{(i)}(n)$ if $i$ does not satisfy (i), (ii) or (iii). If $i$ is in case (i), set $\zeta^{(i)}(n+1) = 0$. For $i$ in case (ii) (resp. (iii)), set $\zeta^{(i)}(n+1) = 1$ (resp. $-1$).

(b) If $\lambda' < \lambda^*$, terminate the algorithm without updating $\zeta$ and set

$$b^{(i)} = b_{\zeta(n)}^{(i)} - \lambda^* u_{\zeta(n)}^{(i)}, \quad \zeta^{(i)}(n) \neq 0$$

and $a_0 = \bar{y} - b^T \bar{x}$ to obtain the final solution.

---

## 7.3   Other Sparsity Estimators

### 7.3.1   LARS estimator

**Algorithm.**   The LARS algorithm can be seen as a simplification of the previous lasso algorithm in which one always adds active variables at each step. We assume as above that input variables are normalized such that $\hat{\sigma}(i) = 1$.

Given a current set $J$ of selected variables, the algorithm will decide either to stop or to add a new variable to $J$ according to a criterion that depends on a parameter $\lambda > 0$. Let $b_{(J)} \in \mathbb{R}^{|J|}$ be the least-square estimator based on variables in $J$

$$b_{(J)} = (\mathcal{X}_c(J)^T \mathcal{X}_c(J))^{-1} \mathcal{X}_c(J)^T \mathcal{Y}_c.$$

Let $b_J \in \mathbb{R}^d$ such that $b_J^{(i)} = b_{(J)}^{(i)}$ for $i \in J$ and 0 otherwise. The covariances between the remaining variables and the residuals are given by

$$r_J^{(i)} = \frac{1}{N} \sum_{k=1}^{N} (y_k - \bar{y} - (x_k - \bar{x})^T b_J)(x_k^{(i)} - \bar{x}^{(i)}), \quad i \notin J.$$

If, for all $i \in J$, $|r_J^{(i)}| \leq \sqrt{\lambda/N}$, the procedure is stopped. Otherwise, one adds to $J$ the variable $i$ such that $|r_J^{(i)}|$ is largest and continues.

**Justification.**   Recall the notation $|b|_0$ for the number of non-zero entries of $b$. Consider the objective function

$$L(b) = |\mathcal{Y}_c - \mathcal{X}_c b|^2 + \lambda |b|_0.$$

Let $J$ be the set currently selected by the algorithm, and $b_J$ defined as above. We consider the problem of adding one non-zero entry to $b$. Fix $i \notin J$, and let $\tilde{b} \in \mathbb{R}^d$ have all coordinates equalt to those of $b_J$ for all except the $i$th one, which is therefore allowed to be non-zero. Then

$$L(\tilde{b}) = \sum_{k=1}^{N} \left( y_k - \overline{y} - \sum_{j \in J} (x_k^{(j)} - \overline{x}) b^{(j)} - (x_k^{(i)} - \overline{x}) \tilde{b}^{(i)} \right)^2 + \lambda |J| + \lambda,$$

so that (using $\hat{\sigma}(i) = 1$)

$$L(\tilde{b}) = L(b_J) - 2N r_J^{(i)} \tilde{b}^{(i)} + N(\tilde{b}^{(i)})^2 + \lambda$$

Now, $L(\tilde{b})$ is an upper-bound for $L(b_{J \cup \{i\}})$, and so is its minimum with respect to $\tilde{b}^{(i)}$. This yields:

$$L(b_{J \cup \{i\}}) \leq L(b_J) - N(r_J^{(i)})^2 + \lambda$$

The LARS algorithm therefore finds the value of $i$ that minimizes this upper-bound, provided that the resulting minimum is less that $L(b_J)$.


**Variant.**   The same argument can be made with $|b|_0$ replaced by $|b|_1$ and one gets

$$L(\tilde{b}) = L(b_J) - 2N r_J^{(i)} \tilde{b}^{(i)} + N(\tilde{b}^{(i)})^2 + \lambda |\tilde{b}^{(i)}|$$

Minimizing this expression with respect to $\tilde{b}^{(i)}$ yields the upper bound:

$$L(b_{J \cup \{i\}}) \leq \begin{cases} L(b_J) - N\left( |r_J^{(i)}| - \dfrac{\lambda}{2N} \right)^2 & \text{if } |r_J^{(i)}| \geq \dfrac{\lambda}{2N} \\[2ex] L(b_J) & \text{if } |r_J^{(i)}| \leq \dfrac{\lambda}{2N} \end{cases}$$


This leads to the following alternate form of LARS. Given a current set $J$ of selected variables, compute

$$r_J^{(i)} = \frac{1}{N} \sum_{k=1}^{N} (y_k - \overline{y} - (x_k - \overline{x})^T b_J)(x_k^{(i)} - \overline{x}^{(i)}), \quad i \notin J.$$

If, for all $i \notin J$, $|r_J^{(i)}| \leq \lambda/2N$, stop the procedure. Otherwise, add to $J$ the variable $i$ such that $|r_J^{(i)}|$ is largest and continue. This form tends to add more variables since the stopping criterion decreases in $1/N$ instead of $1/\sqrt{N}$.

**Why "least angle"?**   Let $\mu_{J,k} = y_k - \overline{y} - (x_k - \overline{x})^T b_J$ denote the residual after regression. The empirical correlation between $\mu$ and $x^{(i)}$ is equal to the cosine of the angle, say $\theta_J^{(i)}$ between $\mu_J \in \mathbb{R}^N$ and $x^{(i)} - \overline{x}$ both considered as vectors in $\mathbb{R}^N$. This cosine is also equal to

$$\cos \theta_J^{(i)} = \frac{\mu_J^T (x^{(i)} - \overline{x}^{(i)})}{|x^{(i)} - \overline{x}^{(i)}||\mu_J|} = \sqrt{N} \frac{r_J^{(i)}}{|\mu_J|}$$

where we have used the fact that $|x^{(i)} - \overline{x}^{(i)}|/\sqrt{N} = \hat{\sigma}(i) = 1$. Since $|\mu_J|$ does not depend on $i$, looking for the largest value of $|r_J^{(i)}|$ is equivalent to looking for the smallest value of $|\theta_J^{(i)}|$, so that we are looking for the unselected variable for which the angle with the current residual is minimal.

### 7.3.2   The Dantzig selector

**Noise-free case.**   Assume that one wants to solve the equation $\mathcal{X}\beta = \mathcal{Y}$ when the dimension, $N$, of $\mathcal{Y}$ is small compared to number of columns, $d$, in $\mathcal{X}$. Since the system is under-determined, one needs additional constraints on $\beta$ and a natural one is to look for sparse solutions, i.e., find solutions with a maximum number of zero coefficients. However, this is numerically challenging, and it is easier to minimize the $\ell^1$ norm of $\beta$ instead (as seen when discussing the lasso, using this norm often provides sparse solutions). In the following, we assume that the empirical variance of each variable is normalized, so that, denoting $\mathcal{X}(i)$ the $i$th column of $\mathcal{X}$, we have $|\mathcal{X}(i)| = 1$.

The Dantzig selector [47] minimizes

$$\sum_{i=1}^d |\beta^{(i)}|$$

subject to the constraint $\mathcal{X}\beta = \mathcal{Y}$. This results in a linear program (therefore easy to implement). More precisely, introducing slack variables, it is indeed equivalent to minimize

$$\sum_{i=1}^d \xi(i) + \sum_{i=1}^d \xi^*(i)$$

subject to constraints $\xi(i) \geq \beta^{(i)}$, $\xi(i)^* \geq -\beta^{(i)}$, $\xi(i) \geq 0$, $\xi^*(i) \geq 0$ and $\mathcal{X}\beta = \mathcal{Y}$.

**Sparsity recovery**   Under some assumptions, this method does recover sparse solutions when they exist. More precisely, let $\hat{\beta}$ be the solution of the linear programming problem above. Assume that there is a set $J^* \subset \{1, \dots, d\}$ such that $\mathcal{X}\beta = \mathcal{Y}$ for some

$\beta \in \mathbb{R}^d$ with $\beta^{(i)} = 0$ if $i \notin J^*$. Conditions under which $\hat{\beta}$ is equal to $\beta$ are provided in Candes and Tao [47] and involve the correlations between pairs of columns of $\mathcal{X}$, and the size of $J$.

That the size of $J^*$ must be a factor is clear, since, for the statement to make sense, there cannot exist two $\beta$'s satisfying $\mathcal{X}\beta = \mathcal{Y}$ and $\beta^{(i)} = 0$ for $i \notin J^*$. Uniqueness is obviously not true if $|J| > N$, because, even if one knew $J$, the condition would be under-constrained for $\beta$. Since the set $J^*$ is not known, and we also want to avoid any other solution associated to a set of same size. So, there cannot exist $\beta$ and $\tilde{\beta}$ respectively vanishing outside of $J^*$ and $\tilde{J}^*$, where $J^*$ and $\tilde{J}^*$ have same cardinality, such that $\mathcal{X}\beta = \mathcal{Y} = \mathcal{X}\tilde{\beta}$. The equation $\mathcal{X}(\beta - \tilde{\beta}) = 0$ would be under-constrained as soon as the number of non-zero coefficients of $\beta - \tilde{\beta}$ is larger than $N$, and since this number can be as large as $|J^*| + |\tilde{J}^*| = 2|J^*|$, we see that one should impose at least $|J^*| \le N/2$.

Given this restriction, another obvious remark is that, if the set $J$ on which $\beta$ does not vanish is known, with $|J|$ small enough, then $\mathcal{X}\beta = \mathcal{Y}$ is over-constrained and any solution is (typically) unique. So the issue really is whether the set $J_\beta$ listing the non-zero indexes of a solution $\beta$ is equal to y $J^*$.

As often, precious insight on the solution of this minimization problem is obtained by considering the dual problem. Introducing Lagrange multipliers $\lambda(i) \ge 0, i = 1, \ldots, d$ for the constraints $\xi(i) - \beta^{(i)} \ge 0$, $\lambda^*(i) \ge 0$, $i = 1, \ldots, d$ for $\xi^*(i) + \beta^{(i)} \ge 0$, $\gamma(i), \gamma^*(i) \ge 0$ for $\xi(i) \ge 0$ and $\xi^*(i) \ge 0$, and $\alpha \in \mathbb{R}^N$ for $\mathcal{X}\beta = \mathcal{Y}$, the Lagrangian is

$$\mathcal{L}(\beta, \xi, \lambda, \lambda^*, \alpha) = (\mathbb{1}_d - \lambda - \gamma)^T \xi + (\mathbb{1}_d - \lambda^* - \gamma^*)^T \xi^* + (\lambda - \lambda^* + \mathcal{X}^T \alpha)^T \beta - \alpha^T \mathcal{Y}.$$

The KKT conditions require $\gamma = \mathbb{1}_d - \lambda$, $\gamma^* = \mathbb{1}_d - \lambda^*$, $\mathcal{X}\alpha = \lambda^* - \lambda$ and the complementary slackness conditions give $(1 - \lambda(i))\xi(i) = (1 - \lambda_i^*)\xi^*(i) = 0$, $\lambda(i)(\beta^{(i)} - \xi(i)) = \lambda^*(i)(\beta^{(i)} + \xi^*(i)) = 0$.

The dual problem requires to minimize $\alpha^T \mathcal{Y}$ subject to the constraints $\mathcal{X}^T \alpha = \lambda^* - \lambda$ and $0 \le \lambda(i), \lambda^*(i) \le 1$. Assume that $(\alpha, \lambda, \lambda^*)$ is a solution of this dual problem. One has the following cases.

(1) If $\lambda(i) \in (0,1)$, then $\xi(i) = \beta^{(i)} - \xi(i) = 0$, which implies $\xi(i) = \beta^{(i)} = 0$, and, as a consequence $(1 - \lambda^*(i))\xi^*(i) = \lambda^*(i)\xi^*(i) = 0$, so that also $\xi^*(i) = 0$ .

(2) Similarly, $\lambda^*(i) \in (0,1)$ implies $\xi(i) = \xi^*(i) = \beta^{(i)} = 0$.

(3) If $\lambda(i) = \lambda^*(i) = 1$, then $\beta^{(i)} - \xi(i) = \beta^{(i)} + \xi(i) = 0$ with $\xi(i), \xi^*(i) \ge 0$, so that also $\xi(i) = \xi^*(i) = \beta^{(i)} = 0$.

(4) If $\lambda(i) = \lambda^*(i) = 0$, then $\xi(i) = \xi^*(i) = 0$ and since $\beta^{(i)} \le \xi(i)$ and $\beta^{(i)} \le -\xi^*(i)$, we get $\beta^{(i)} = 0$.

(5) The only remaining situation, in which $\beta^{(i)}$ can be non-zero, is when $\lambda(i) = 1 - \lambda^*(i) \in \{0, 1\}$, or, equivalently, when $|\lambda(i) - \lambda^*(i)| = 1$.

This discussion allows one to reconstruct the set $J_\beta$ associated with the primal problem given the solution of the dual problem. Note that $|\lambda(i) - \lambda^*(i)| = |\alpha^T \mathcal{X}(i)|$, so that the set of indexes with $|\lambda(i) - \lambda^*(i)| = 1$ is also

$$I_\alpha \overset{\Delta}{=} \left\{ i : |\alpha^T \mathcal{X}(i)| = 1 \right\}.$$

One has

$$\alpha^T \mathcal{Y} = \alpha^T \mathcal{X}\beta = \sum_{i=1}^{d} \beta^{(i)} \alpha^T \mathcal{X}(i) \le \sum_{i \in J_\beta} |\beta^{(i)}| |\alpha^T \mathcal{X}(i)| \le \sum_{i \in J_\beta} |\beta^{(i)}|.$$

The upper-bound is achieved when $\alpha^T \mathcal{X}(i) = \text{sign}(\beta^{(i)})$ for $i \in J_\beta$. So, if a vector $\alpha$ can be found such that

(i) $\alpha^T \mathcal{X}(i) = \text{sign}(\beta^{(i)})$ for $i \in J^*$,

(ii) $|\alpha^T \mathcal{X}(j)| < 1$ for $j \notin J^*$,

then it is a solution of the dual problem with $J_\alpha = J^*$.

Let $s_J = (s^{(j)}, j \in J)$ be defined by $s^{(j)} = \text{sign}(\beta^{(j)})$. One can always decompose $\alpha \in \mathbb{R}^N$ in the form

$$\alpha = \mathcal{X}_{J^*} \rho + w$$

where $\rho \in \mathbb{R}^{|J^*|}$ and $w \in \mathbb{R}^N$ is perpendicular to the columns of $\mathcal{X}_{J^*}$. From $\mathcal{X}_{J^*}^T \alpha = s_J$, we get

$$\rho = (\mathcal{X}_{J^*}^T \mathcal{X}_{J^*})^{-1} s_{J^*}.$$

Letting $\alpha_{J^*}$ be the solution with $w = 0$, the question is therefore whether one can find $w$ such that

$$\begin{cases} w^T \mathcal{X}(j) = 0, & j \in J^* \\ |\alpha_J^T \mathcal{X}(k) + w^T \mathcal{X}(k)| < 1, & k \notin J^* \end{cases}$$

Denote for short $\Sigma_{JJ'} = \mathcal{X}_J^T \mathcal{X}_{J'}$. One can show that such a solution exists when the matrices $\Sigma_{JJ}$ are close to the identity as soon as $|J|$ is small enough [47]. More precisely, denote, for $q \le d$

$$\delta(q) = \max_{|J| \le q} \max(\|\Sigma_{JJ}\|, \|\Sigma_{JJ}^{-1}\|^{-1}) - 1,$$

in which one uses the operator norm on matrices, and

$$\theta(q, q') = \max\left\{ z^T \Sigma_{TT'} z' : |J|, |J'| \le q, J \cap J' = \emptyset, |z| = |z'| = 1 \right\}.$$

Then, the following proposition is true.

**Proposition 7.7 (Candes-Tao)** *Let $q = |J^*|$ and $s = (s^{(j)}, j \in J^*) \in \mathbb{R}^q$. Assume that $\delta(2q) + \theta(q, 2q) < 1$. Then there exists $\alpha \in \mathbb{R}^N$ such that $\alpha^T \mathcal{X}(j) = s^{(j)}$ for $j \in J^*$ and*

$$|\alpha^T \mathcal{X}(j)| \le \frac{\theta(q, q)}{1 - \delta(2q) - \theta(q, 2q)} \text{ if } j \notin J^*.$$

So $\alpha$ has the desired property as soon as $\delta(2q) + \theta(q, q) + \theta(q, 2q) \le 1$. to control subsets of variables of size less than $3q$ to obtain the conclusion, which is important, of course, when $q$ is small compared to $d$.

**Noisy case**   Consider now the noisy case. We here again introduce quantities that were pivotal for the lasso and LARS estimators, namely, the covariances between the variables and the residual error. So, we define, for a given $\beta$

$$r_\beta^{(i)} = \mathcal{X}(i)^T (\mathcal{Y} - \mathcal{X}\beta)$$

which depends linearly on $\beta$. Then, the *Dantzig selector* is defined by the linear program: Minimize:

$$\sum_{j=1}^{d} |\beta^{(j)}|$$

subject to the constraint:

$$\max_{j=1,\dots,d} |r_\beta^{(j)}| \le C.$$

The explicit expression of this problem as a linear program is obtained as before by introducing slack variables $\xi(j), \xi^*(j), j = 1, \dots, d$ and minimizing

$$\sum_{j=1}^{d} \xi(j) + \sum_{j=1}^{d} \xi^*(j)$$

with constraints $\xi(j), \xi^*(j) \ge 0, \xi \ge \beta, \xi^* \ge -\beta, \max_{j=1,\dots,d} |r_\beta^{(j)}| \le C.$

Similar to the noise-free case, the Dantzig selector can identify sparse solutions (up to a small error) if the columns of $\mathcal{X}$ are nearly orthogonal, with the same type of conditions [48]. Interestingly enough, the accuracy of this algorithm can be proved to be comparable to that of the lasso in the presence of a sparse solution [30].

## 7.4 Support Vector Machines for regression

### 7.4.1 Linear SVM

**Problem formulation** We start by discussing support vector machines (SVM) [196, 197] with $\mathcal{R}_X = \mathbb{R}^d$ equipped with the standard inner product (generally referred to as linear SVM) and will extend the theory to kernel methods in the next section. SVMs solve a linear regression problem, but replace the least-squares loss function by $(y, y') \mapsto V(y - y')$ with

$$V(t) = \begin{cases} 0 & \text{if } |t| < \epsilon \\ |t| - \epsilon & \text{if } |t| \geq \epsilon \end{cases} \tag{7.10}$$



Figure 7.1: The function $V$ defining the SVM risk function.

A plot of the function $V$ is provided in fig. 7.1. This function is an example of what is often called a robust loss function. The quadratic error used in linear regression had the advantage of providing closed form expressions for the solution, but is quite sensitive to outliers. For robustness, it is preferable to use loss functions that, like $V$, increase at most linearly at infinity. One sometimes choose them as smooth convex functions, for example $V(t) = (1 - \cos \gamma t)/(1 - \cos \gamma \epsilon)$ for $|t| < \epsilon$ and $f(t) = |t|$ for $t \geq \epsilon$, where $\gamma$ is chosen so that $\gamma \epsilon \sin \gamma \epsilon / (1 - \cos \gamma \epsilon) = 1$. In such a case, minimizing

$$F(\beta) = \sum_{k=1}^{N} V(y_k - a_0 - x_k^T b)$$

can be done using gradient descent methods. Using $V$ in (7.10) will require a little more work, as we see now.

The SVM regression problem is generally formulated as the minimization of

$$\sum_{k=1}^{N} V(y_k - a_0 - x_k^T b) + \lambda |b|^2,$$

and we will study a slightly more general problem, minimizing

$$F(a_0, b) = \sum_{k=1}^{N} V(y_k - a_0 - x_k^T b) + \lambda b^T \Delta b,$$

where $\Delta$ is a symmetric positive-definite matrix. This objective function exhibits the following features:

- A penalty on the coefficients of $b$, similar to ridge regression.

- A linear penalty (instead of quadratic) for large errors in the prediction.

- An $\epsilon$-tolerance for small errors, often referred to as *the margin* of the regression SVM.

We now describe the various steps in the analysis and reduction of the problem. They will lead to simple minimization algorithms, and possible extensions to non-linear problems.

**Reduction to a quadratic programming problem**   Introduce slack variables $\xi_k, \xi_k^*, k = 1, \dots, N$. The original problem is equivalent to the minimization, with respect to $(a_0, b, \xi, \xi^*)$, of

$$\sum_{k=1}^{N} (\xi_k + \xi_k^*) + \lambda b^T \Delta b$$

under the constraints:

$$\begin{cases} \xi_k, \xi_k^* \geq 0 \\ \xi_k - y_k + a_0 + x_k^T b + \epsilon \geq 0 \\ \xi_k^* + y_k - a_0 - x_k^T b + \epsilon \geq 0 \end{cases} \tag{7.11}$$

The simple proof of this equivalence, which results in a quadratic programming problem, is left to the reader. As often, one gains additional insight by studying the dual problem.

**Dual problem**   Introduce $4N$ non-negative Lagrange multipliers for the $4N$ constraints in the problem, namely, $\eta_k, \eta_k^* \geq 0$ for the positivity constraints, and $\alpha_k, \alpha_k^* \geq$

0 for the last two in (7.11). The resulting Lagrangian is

$$\mathcal{L}(a_0, b, \xi, \xi^*, \alpha, \alpha^*, \eta, \eta^*) = \sum_{k=1}^{N} (\xi_k + \xi_k^*) + \lambda b^T \Delta b - \sum_{k=1}^{N} (\eta_k \xi_k + \eta_k^* \xi_k^*)$$
$$- \sum_{k=1}^{N} \alpha_k (\xi_k - y_k + a_0 + x_k^T b + \epsilon) - \sum_{k=1}^{N} \alpha_k^* (\xi_k^* + y_k - a_0 - x_k^T b + \epsilon).$$

In this formulation, $(a_0, b, \xi, \xi^*)$ are the primal variables, and $\alpha, \alpha^*, \eta, \eta^*$ the dual variables.

The KKT conditions are provided by the system:

$$\begin{cases} \sum_{k=1}^{N} (\alpha_k - \alpha_k^*) = 0 \\ 2\lambda \Delta b - \sum_{k=1}^{N} (\alpha_k - \alpha_k^*) x_k = 0 \\ 1 - \eta_k - \alpha_k = 0 \\ 1 - \eta_k^* - \alpha_k^* = 0 \\ \alpha_k (\epsilon + \xi_k - y_k + a_0 + x_k^T b) = 0 \\ \alpha_k^* (\epsilon + \xi_k^* + y_k - a_0 - x_k^T b) = 0 \\ \eta_k \xi_k = \eta_k^* \xi_k^* = 0 \end{cases} \tag{7.12}$$

The first four equations are the derivatives of the Lagrangian with respect to $a_0, b, \xi_k, \xi_k^*$ in this order and the last three are the complementary slackness conditions.

The dual problem maximizes the function

$$L^*(\alpha, \alpha^*, \eta, \eta^*) = \inf_{\beta, \xi, \xi^*} \mathcal{L}.$$

under the previous positivity constraints. Since the Lagrangian is linear in $a_0$, $\xi_k$ and $\xi_k^*$, its minimum is $-\infty$ unless the coefficients vanish. The linear terms must therefore vanish for $L^*$ to be finite. With these conditions plus the fact that $\partial_b L = 0$, we retrieve the first four equations of system (7.12). Using $\eta_k = 1 - \alpha_k$, $\eta_k^* = 1 - \alpha_k^*$ and

$$b = \frac{1}{2\lambda} \sum_{k=1}^{N} (\alpha_k - \alpha_k^*) \Delta^{-1} x_k \tag{7.13}$$

one can express $L^*$ uniquely as a function of $\alpha, \alpha^*$, yielding

$$L^*(\alpha, \alpha^*) = -\frac{1}{4\lambda} \sum_{k,l=1}^{N} (\alpha_k - \alpha_k^*)(\alpha_l - \alpha_l^*) x_k^T \Delta^{-1} x_l - \epsilon \sum_{k=1}^{N} (\alpha_k + \alpha_k^*) + \sum_{k=1}^{N} (\alpha_k - \alpha_k^*) y_k.$$

This quantity must be maximized subject to the constraints $0 \leq \alpha_k, \alpha_k^* \leq 1$ and $\sum_{k=1}^{N}(\alpha_k - \alpha_k^*) = 0$. This still is a quadratic programming problem, but it now has nice additional features and interpretations.

**Step 3: Analysis of the dual problem**    The dual problem only depends on the $x_k$'s through the matrix with coefficients $x_k^T \Delta^{-1} x_l$, which is the Gram matrix of $x_1, \ldots, x_N$ for the inner product associated with $\Delta^{-1}$. This property will lead to the the kernel version of SVMs discussed in the next section. The obtained predictor can also be expressed as a function of these products, since

$$y = a_0 + x^T b = a_0 + \frac{1}{2\lambda} \sum_{k=1}^{N} (\alpha_k - \alpha_k^*)(x_k^T \Delta^{-1} x).$$

Moreover, the dimension of the dual problem is $2N$, which allows the method to be used in large (possibly infinite) dimensions with a bounded cost.

We now analyze the solutions $\alpha, \alpha^*$ of the dual problem. The complementary slackness conditions reduce to:

$$\begin{cases} \alpha_k(\epsilon + \xi_k - y_k + a_0 + x_k^T b) = 0 \\ \alpha_k^*(\epsilon + \xi_k^* + y_k - a_0 - x_k^T b) = 0 \\ (1 - \alpha_k)\xi_k = (1 - \alpha_k^*)\xi_k^* = 0 \end{cases} \quad (7.14)$$

These conditions have the following consequences, based on the prediction error made for each training sample.

(i) First consider indexes $k$ such that the error is strictly within the tolerance margin $\epsilon$: $|y_k - a_0 - x_k^T b| < \epsilon$. Then the terms between parentheses in first two equations of (7.14) are strictly positive, which implies that $\alpha_k = \alpha_k^* = 0$. The last two equations in (7.14) then imply $\xi_k = \xi_k^* = 0$.

(ii) Consider now the case when the prediction is strictly less accurate than the tolerance margin. Assume that $y_k - a_0 - x_k^T b > \epsilon$. The second and third equations in (7.14) imply that $\alpha_k^* = \xi_k^* = 0$. The assumption also implies that

$$\xi_k = y_k - a_0 - x_k^T b - \epsilon > 0$$

and $\alpha_k = 1$. The case $y_k - a_0 - x_k^T b < -\epsilon$ is symmetric and provides $\alpha_k = \xi_k = 0$, $\xi_k^* > 0$ and $\alpha_k^* = 1$.

(iii) Finally, consider samples for which the prediction error is exactly at the tolerance margin. If $y_k - a_0 - x_k^T b = \epsilon$, we have $\alpha_k^* = \xi_k = \xi_k^* = 0$. The fact that $\alpha_k^* = \xi_k^* = 0$ is clear. To prove that $\xi_k = 0$, we note that would have otherwise $\xi_k - y_k + a_0 + x_k^T b + \epsilon > 0$, which would imply that $\alpha_k = 0$ and we reach a contradiction with $(1 - \alpha_k)\xi_k = 0$. Similarly, $y_k - a_0 - x_k^T b = -\epsilon$ implies that $\alpha_k = \xi_k = \xi_k^* = 0$.

The points for which $|y_k - a_0 - x_k^T b| = \epsilon$ are called *support vectors*.

One important information deriving from this discussion is that the variables $(\alpha_k, \alpha_k^*)$ have prescribed values as long as the error $y_k - a_0 - x_k^T b$ is not exactly $\epsilon$ in absolute value: $(1, 0)$ if the error is larger than $\epsilon$, $(0, 0)$ if it is strictly between $-\epsilon$ and $\epsilon$ and $(0, 1)$ if it is less than $-\epsilon$. Also in all cases, at least one of $\alpha_k$ and $\alpha_k^*$ must vanish. Only in the case of support vectors does the previous discussion fail to provide a value for one of these variables.

Now, we want to reverse the discussion and assume that the dual problem is solved to see how the variables $a_0$ and $b$ of the primal problem can be retrieved. For $b$, this is easy, thanks to (7.13). For $a_0$ a direct computation can be made if a support vector is identified, either because $0 < \alpha_k < 1$, which implies that $a_0 = y_k - x_k^T b - \epsilon$, or because $0 < \alpha_k^* < 1$, which yields $a_0 = y_k - x_k^T b + \epsilon$.

If no support vector can be identified, $a_0$ is not uniquely determined (note that the objective function is not strictly convex in $a_0$). However, the coefficients $\alpha_k, \alpha_k^*$ provide some information on this intercept, in the form of inequalities. More precisely, let $J^+ = \{k : \alpha_k = 1\}$, $J^- = \{k : \alpha_k^* = 1\}$ and $J_0 = \{k : \alpha_k = \alpha_k^* = 0\}$. Then $k \in J^+$ implies that $y_k - a_0 - b^T x_k \geq \epsilon$, so that $a_0 \leq y_k - b^T x_k - \epsilon$. Similarly, $k \in J^-$ implies that $a_0 \geq y_k - b^T x_k + \epsilon$. Finally, $k \in J_0$ implies that $a_0 \geq y_k - b^T x_k - \epsilon$ and $a_0 \leq y_k - b^T x_k + \epsilon$. As a consequence, one can take $a_0$ to be any point in the interval $[a_0^-, a_0^+]$, where

$$a_0^- = \max\left(\max_{k \in J^-}(y_k - x_k^T b + \epsilon), \max_{k \in J_0}(y_k - x_k^T b - \epsilon)\right)$$

$$a_0^+ = \min\left(\min_{k \in J^+}(y_k - x_k^T b - \epsilon), \min_{k \in J_0}(y_k - x_k^T b + \epsilon)\right).$$

### 7.4.2 The kernel trick and SVMs

Returning to our feature space notation, let $X$ take values in $\mathcal{R}_X$ and $h : \mathcal{R}_X \to H$ be a feature function with values in an inner-product space $H$ with associated kernel $K$. SVMs in feature space must minimize, with $a_0 \in \mathbb{R}$ and $b \in H$

$$F(a_0, b) = \sum_{k=1}^{N} V\left(y_k - a_0 - \langle h(x_k), b \rangle_H\right) + \lambda \|b\|_H^2.$$

Letting as before $V = \text{span}(h(x_1), \ldots, h(x_N))$, the same argument as that made for ridge regression works, namely that the first term in $F$ is unchanged if $b$ is replaced by $\pi_V(b)$ and the second one is strictly reduced unless $b \in V$, leading to a finite-dimensional formulation in which

$$b = \sum_{k=1}^{N} c_k h(x_k)$$

and one minimizes

$$F(a_0, c) = \sum_{k=1}^{N} V\left(y_k - a_0 - \sum_{l=1}^{N} K(x_k, x_l)c_l\right) + \lambda \sum_{k,l=1}^{N} K(x_k, x_l)c_k c_l.$$

This function has the same form as the one studied in the linear case with $b$ replaced by $c \in \mathbb{R}^N$, $x_k$ replaced by the vector with coefficients $K(x_k, x_l), l = 1, \ldots, N$, that we will denote $\mathcal{K}^{(k)}$ and $\Delta = \mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$. Note that $\mathcal{K}^{(k)}$ is the $k$th column of $\mathcal{K}$, so that

$$\left(\mathcal{K}^{(k)}\right)^T \mathcal{K}^{-1} \mathcal{K}^{(l)} = K(x_k, x_l).$$

Using this, we find that the dual problem requires to maximize

$$L^*(\alpha, \alpha^*) = -\frac{1}{4\lambda} \sum_{k,l=1}^{N} (\alpha_k - \alpha_k^*)(\alpha_l - \alpha_l^*)K(x_k, x_l) - \epsilon \sum_{k=1}^{N} (\alpha_k + \alpha_k^*) + \sum_{k=1}^{N} (\alpha_k - \alpha_k^*)y_k.$$

with

$$\begin{cases} 0 \leq \alpha_k \leq 1 \\ 0 \leq \alpha_k^* \leq 1 \\ \displaystyle\sum_{k=1}^{N} (\alpha_k - \alpha_k^*) = 0 \end{cases}$$

The associated vector $c$ satisfies

$$2\lambda c = \sum_{k=1}^{N} (\alpha_k - \alpha_k^*)\mathcal{K}^{-1}\mathcal{K}^{(k)} = \alpha - \alpha^*.$$

and the regression function is

$$f(x) = a_0 + \langle b, h(x) \rangle_H = a_0 + \frac{1}{2\lambda} \sum_{k=1}^{N} (\alpha_k - \alpha_k^*)K(x, x_k).$$

Finally, the discussions on the values of $\alpha, \alpha^*$ and on the computation of $a_0$ remain unchanged.

# Chapter 8

# Models for linear classification

In this chapter, $Y$ is categorical and takes values in the finite set $\mathcal{R}_Y = \{g_1,\ldots,g_q\}$. The goal is to predict this class variable from the input $X$, taking values in a set $\mathcal{R}_X$. Using the same progression as in the regression case, we will first discuss basic linear methods, for which $\mathcal{R}_X = \mathbb{R}^d$ before extending them, whenever possible, to kernel methods, for which $\mathcal{R}_X$ can be arbitrary as soon as a feature space representation is available.

Classifiers will be based on a training set $T = ((x_1,y_1),\ldots,(x_N,y_N))$ with $x_k \in \mathcal{R}_X$ and $y_k \in \mathcal{R}_Y$ for $k = 1,\ldots,N$. For $g \in \mathcal{R}_Y$, we will also let $N_g$ denote the number of samples in the training set such that $y_k = g$, i.e.,

$$N_g = \left|\{k : y_k = g\}\right| = \sum_{k=1}^{N} \mathbf{1}_{y_k=g}.$$

## 8.1 Logistic regression

### 8.1.1 General Framework

Logistic regression uses the fact that, in order to apply Bayes's rule, only the conditional distribution of the class variables $Y$ given $X$ is needed, and trains a parametric model of this distribution. More precisely, if one denotes by $p(g|x)$ the probability that $Y = g$ conditional to $X = x$, logistic regression assumes that, for some parameters $(a_0(g),b(g),g \in \mathcal{R}_Y)$ with $a_0(g) \in \mathbb{R}$ and $b(g) \in \mathbb{R}^d$, one has $p = p_{a_0,b}$ with

$$\log p_{a_0,b}(g \mid x) = a_0(g) + x^T b(g) - \log(C(a_0,b,x)),$$

where $C(a_0,b,x) = \sum_{g \in \mathcal{R}_Y} \exp(a_0(g) + x^T b(g))$.

Introduce the functions, defined over mappings $\mu : \mathcal{R}_Y \to \mathbb{R}$ (which can be identified with vectors in $\mathbb{R}^q$)

$$F_g(\mu) = \mu(g) - \log \sum_{g' \in \mathcal{R}_Y} e^{\mu(g')}. \tag{8.1}$$

With this notation, letting $\beta(g) = \begin{pmatrix} a_0(g) \\ b(g) \end{pmatrix}$ and $\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$, one has $\log p_\beta(g|x) = F_g(\beta^T \tilde{x})$, where $\beta^T \tilde{x}$ is the function $(g' \mapsto \beta(g')^T \tilde{x})$.

For any constant function $(g \mapsto \mu_0 \in \mathbb{R})$ one has

$$F_g(\mu + \mu_0) = \mu(g) + \mu_0 - \log \sum_{g' \in \mathcal{R}_Y} e^{\mu(g') + \mu_0} = \mu(g) + \mu_0 - \mu_0 - \log \sum_{g' \in \mathcal{R}_Y} e^{\mu(g')} = F_g(\mu).$$

As a consequence, if one replaces, for all $g$, $\beta(g)$ by $\tilde{\beta}(g) = \beta(g) + \gamma$, with $\gamma \in \mathbb{R}^{d+1}$, then $\tilde{\beta}^T \tilde{x} = \beta^T \tilde{x} + \gamma^T \tilde{x}$ and

$$\log p_{\tilde{\beta}}(g \mid x) = \log p_\beta(g \mid x).$$

This shows that the model is over-parametrized. One therefore needs a $(d+1)$-dimensional constraint to ensure uniqueness, and we will enforce a linear constraint in the form

$$\sum_{g \in \mathcal{R}_Y} \rho_g \beta(g) = c$$

with $\sum_g \rho_g \neq 0$.

### 8.1.2   Conditional log-likelihood

The conditional log-likelihood computed from the training set is:

$$\ell(\beta) = \sum_{k=1}^N \log p_\beta(y_k \mid x_k).$$

Logistic regression computes a maximizer $\hat{\beta}$ of this log-likelihood. The classification rule given a new input $x$ then chooses the class $g$ for which $p_{\hat{\beta}}(g \mid x)$ is largest, or, equivalently, the class $g$ that maximizes $\tilde{x}^T \beta(g)$.

**Proposition 8.1** *Let $\tilde{m}_g = \sum_{k:y_g=k} x_k / N_g$. The conditional log-likelihood $\ell$ is concave with first derivative*

$$\partial_{\beta(g)} \ell = N_g \tilde{m}_g^T - \sum_{k=1}^N \tilde{x}_k^T p_\beta(g|x_k) \tag{8.2}$$

*and negative semi-definite second derivative*

$$\partial_{\beta(g)}\partial_{\beta(g')}\ell = -\mathbf{1}_{[g=g']}\sum_{k=1}^{N}\tilde{x}_k\tilde{x}_k^T p_\beta(g|x_k) + \sum_{k=1}^{N}\tilde{x}_k\tilde{x}_k^T p_\beta(g|x_k)p_\beta(g'|x_k). \tag{8.3}$$

**Remark 8.2** In this discussion, we consider $\ell$ as a function defined over collections $(\beta(g), g \in \mathcal{R}_Y)$, or, if one prefers, on the $q(d+1)$-dimensional linear space, $\mathcal{F}$, of functions $\beta : \mathcal{R}_Y \to \mathbb{R}^{q+1}$. With this in mind, the differential $d\ell(\beta)$ is a linear form from $\mathcal{F}$ to $\mathbb{R}$, therefore associating to any family $u = (u(g), g \in \mathcal{R}_Y)$, the expression

$$d\ell(\beta)u = \sum_{g \in \mathcal{R}_Y} \partial_{\beta(g)}\ell \, u(g).$$

Similarly, the second derivative is the bilinear form

$$d^2\ell(\beta)(u, u') = \sum_{g,g' \in \mathcal{R}_Y} u(g)^T \partial_{\beta(g)}\partial_{\beta(g')}\ell \, u(g').$$

The last statement in the proposition expresses the fact that $d^2\ell(\beta)(u, u) \leq 0$ for all $u \in \mathcal{F}$. ♦

PROOF First consider the function $F_g$ in (8.1), so that

$$\ell(\beta) = \sum_{k=1}^{N} F_{y_k}(\beta^T \tilde{x}_k).$$

We have, for $\zeta : \mathcal{R}_Y \to \mathbb{R}$,

$$dF_g(\mu)\zeta = \zeta(g) - \frac{\sum_{g' \in \mathcal{R}_Y} e^{\mu(g')}\zeta(g')}{\sum_{g' \in \mathcal{R}_Y} e^{\mu(g')}}$$

as can be easily computed by evaluating the derivative of $F(\mu + \epsilon u)$ at $\epsilon = 0$. Introducing the notation

$$q_\mu(g) = \frac{e^{\mu(g)}}{\sum_{g \in \mathcal{R}_Y} e^{\mu(g)}}$$

and

$$\langle \zeta \rangle_\mu = \sum_{g \in \mathcal{R}_Y} \zeta(g)q_\mu(g),$$

we have $dF_g(\mu)\zeta = \zeta(g) - \langle \zeta \rangle_\mu$. Evaluating the derivative of $dF_g(\mu + \epsilon u')(\zeta)$ at $\epsilon = 0$, one gets (the computation being left to the reader)

$$d^2F_g(\mu)(\zeta, \zeta') = -\langle \zeta\zeta' \rangle_\mu + \langle \zeta \rangle_\mu \langle \zeta' \rangle_\mu. \tag{8.4}$$

Note that $-d^2 F_g(\mu)(\zeta, \zeta)$ is the variance of $\zeta$ for the probability mass function $q_\mu$ and is therefore non-negative (so that $F_\mu$ is concave). This immediately shows that $\ell$ is concave as a sum of concave functions.

Using the chain rule, we have, for $u : \mathcal{R}_Y \to \mathbb{R}^q$,

$$d\ell(\beta)u = \sum_{k=1}^{N} dF_{y_k}(\beta^T \tilde{x}_k) \tilde{x}_k^T u(\cdot) = \sum_{k=1}^{N} \tilde{x}_k^T u(y_k) - \sum_{k=1}^{N} \langle \tilde{x}_k^T u(\cdot) \rangle_{\beta^T \tilde{x}_k}.$$

Reordering the first sum in the right-hand side according to the values of $y_k$ gives

$$\sum_{k=1}^{N} u(y_k)^T \tilde{x}_k = \sum_{g \in \mathcal{R}_Y} N_g u(g)^T \tilde{m}_g.$$

Noting that $q_{\beta^T \tilde{x}} = p_\beta(\cdot | x)$, we find

$$d\ell(\beta)(u) = \sum_{g \in \mathcal{R}_Y} N_g \tilde{m}_g^T u(g) - \sum_{g \in \mathcal{R}_Y} \tilde{x}_k^T u(g) p_\beta(g | x_k),$$

yielding (8.2). Applying the chain rule again, we have

$$d^2\ell(\beta)(u, u') = \sum_{k=1}^{N} d^2 F_{y_k}(\beta^T \tilde{x}_k)(\tilde{x}_k^T u(\cdot), \tilde{x}_k^T u'(\cdot)) \tag{8.5}$$

with

$$\begin{aligned}
d^2 F_{y_k}(\beta^T \tilde{x}_k)(\tilde{x}_k^T u(\cdot), \tilde{x}_k^T u'(\cdot)) &= -\langle u(\cdot)^T \tilde{x}_k \tilde{x}_k^T u'(\cdot) \rangle_{\beta^T \tilde{x}_k} + \langle \tilde{x}_k^T u(\cdot) \rangle_{\beta^T \tilde{x}_k} \langle \tilde{x}_k^T u'(\cdot) \rangle_{\beta^T \tilde{x}_k} \\
&= -\sum_{g \in \mathcal{R}_Y} u(g)^T \tilde{x}_k \tilde{x}_k^T u'(g) p_\beta(g | x_k) \\
&\quad + \sum_{g, g' \in \mathcal{R}_Y} u(g)^T \tilde{x}_k \tilde{x}_k^T u'(g') p_\beta(g | x_k) p_\beta(g' | x_k)
\end{aligned}$$

from which (8.3) follows. ∎

**Remark 8.3** From $F_g(\mu + \mu_0) = F_g(\mu)$ when $\mu_0$ is constant on $\mathcal{R}_Y$, one deduces (taking the derivative at $\mu_0 = 0$) that $dF_g(\mu)\mathbf{1} = 0$ for all $\mu$, where $\mathbf{1}$ denotes the constant function equal to 1 on $\mathcal{R}_Y$. For $h \in \mathbb{R}^{d+1}$, let $\mathbb{c}_h$ denote the constant function $\mathbb{c}_h(g) = h$, $g \in \mathcal{R}_Y$. We have

$$d\ell(\beta)\mathbb{c}_h = \sum_{k=1}^{N} dF_{y_k}(\beta^T \tilde{x}_k) \tilde{x}_k^T \mathbb{c}_h = \sum_{k=1}^{N} dF_{y_k}(\beta^T \tilde{x}_k)(\tilde{x}_k^T h \mathbf{1}) = \sum_{k=1}^{N} (\tilde{x}_k^T h) dF_{y_k}(\beta^T \tilde{x}_k)\mathbf{1} = 0.$$

Taking one extra derivative we see that

$$d\ell(\beta)(\mathbb{c}_h, u) = 0$$

for all functions $u : \mathcal{R}_Y \to \mathbb{R}^q$. ◆

We now discuss whether there are other elements in the null space of the second derivative of $\ell$. We will use notation introduced in the proof of proposition 8.1. From (8.4), we have $d^2 F_g(\mu)(\zeta, \zeta) = 0$ if and only if the variance of $\zeta$ for $q_\mu$ vanishes, which, since $q_\mu > 0$, is equivalent to $\zeta$ being constant. So, the null space of $d^2 F_g(\mu)$ is one-dimensional, and composed of scalar multiples of $\mathbf{1}$. Using (8.5), we see that $d^2\ell(u,u) = 0$ if and only if , for all $k = 1,\dots,N$, $(g \mapsto \tilde{x}_k^T u(g))$ is a constant function.

Assume that this is true. Then, letting $\bar{u} = \frac{1}{q}\sum_{g\in\mathcal{R}_Y} u(g)$, one has, for all $g \in \mathcal{R}_Y$ and $k = 1,\dots,N$,

$$\tilde{x}_k^T u(g) = \tilde{x}_k^T \bar{u}$$

so that $u(g) - \bar{u}$ is in the null space of the matrix $\mathcal{X}$. This leads to the following proposition.

**Proposition 8.4** *Assume that $\mathcal{X}$ has rank $d + 1$. Then the null space of $d^2\ell(\beta)$ is the set of all vectors $u = \mathbb{c}_h$ for $h \in \mathbb{R}^{d+1}$. In particular, for any $c \in \mathbb{R}^{d+1}$, the function $\ell$ restricted to the space*

$$M = \left\{\beta : \sum_{g\in\mathcal{R}_Y} \rho_g\beta(g) = c\right\}$$

*is strictly concave as soon as the scalar coefficients $(\rho_g, g \in \mathcal{R}_Y)$ are such that $\sum_{g\in\mathcal{R}_Y} \rho_g \neq 0$.*

PROOF From the discussion before the proposition, $u \in \mathrm{Null}(d^2\ell)$ implies that $\mathcal{X}(u(g) - \bar{u}) = 0$ for all $g$, and since we assume that $\mathcal{X}$ has rand $d + 1$, this requires that $u(g) = \bar{u}$ for all $g$, i.e., $u = \mathbb{c}_{\bar{u}}$. This proves the first point.

If one restricts $\ell$ to $M$, then we must restrict $d^2\ell(\beta)$ to those $u$'s such that $\sum_{g\in\mathcal{R}_Y} \rho_g u(g) = 0$. But if $d^2\ell(\beta)(u,u) = 0$ for such an $u$, then $u = \mathbb{c}_{\bar{u}}$ and

$$\sum_{g\in\mathcal{R}_Y} \rho_g u(g) = \left(\sum_{g\in\mathcal{R}_Y} \rho_g\right)\bar{u}.$$

Since we assume that $\sum_{g\in\mathcal{R}_Y} \rho_g \neq 0$, this requires $\bar{u} = 0$, and therefore $u = 0$.

This shows that the second derivative of the restriction of $\ell$ to $M$ is negative definite, so this restriction is strictly concave. ∎

### 8.1.3 Training algorithm

Given that we have expressed the first and second derivatives of $\ell$ in closed form[1], we can use Newton-Raphson gradient ascent to maximize $\ell$ over the affine space:

$$M = \left\{\beta : \sum_{g\in\mathcal{R}_Y} \rho_g\beta(g) = c\right\}$$

---

[1]Their computation is feasible unless $N$ is very large, and the matrix inversion in Newton's iteration also requires $d$ to be not too large.

with $\sum_{g \in \mathcal{R}_Y} \rho_g \neq 0$. We assume in the following that the matrix $\mathcal{X}$ has rand $d + 1$ so that proposition 8.4 applies. Since the constraint is affine, it is easy to express one of the parameters $\beta(g)$ as a function of the others and solve the strictly concave problem as a function of the remaining variables. It is not much harder, and arguably more elegant to solve the problem without breaking its symmetry with respect to the class indexes, as described below.

Let

$$M_0 = \left\{ \beta : \sum_{g \in \mathcal{R}_Y} \rho_g \beta(g) = 0 \right\}.$$

We still have the second order expansion

$$\ell(\beta + u) = \ell(\beta) + d\ell(\beta)u + \frac{1}{2} d^2 \ell(\beta)(u, u) + o(|u|^2)$$

and we consider the maximization of the first three terms, simply restricting to vectors $u \in M_0$. To allow for matrix computation, we use our ordering $\mathcal{R}_Y = (g_1, \ldots, g_q)$ and identify $a$ with the column vector

$$\begin{pmatrix} u(g_1) \\ \vdots \\ u(g_q) \end{pmatrix} \in \mathbb{R}^{q(d+1)}$$

Similarly, we let

$$\nabla \ell(\beta) = \begin{pmatrix} \partial_{\beta(g_1)} \ell \\ \vdots \\ \partial_{\beta(g_q)} \ell \end{pmatrix}$$

and let $\nabla^2(\ell)(\beta)$ be the block matrix with $i, j$ block given by $\partial_{\beta(g_i)} \partial_{\beta(g_j)} \ell(\beta)$. We let $\hat{\rho}$ be the $(d + 1) \times q(d + 1)$ row block matrix

$$\begin{pmatrix} \rho(g_1) \mathrm{Id}_{\mathbb{R}^{d+1}} & \cdots & \rho(g_q) \mathrm{Id}_{\mathbb{R}^{d+1}} \end{pmatrix}$$

so that $u \in M_0$ is just $\hat{\rho} u = 0$ in vector notation. Given this we have

$$\ell(\beta + u) = \ell(\beta) + \nabla \ell(\beta)^T u + \frac{1}{2} u^T \nabla^2(\ell)(\beta) u + o(|u|^2).$$

The maximum of $\ell(\beta) + u^T \nabla \ell(\beta) + \frac{1}{2} u^T \nabla^2(\ell)(\beta) u$ subject to $\hat{\rho} u = 0$ is a stationary point of the Lagrangian

$$L = \ell(\beta) + u^T \nabla \ell(\beta) + \frac{1}{2} u^T \nabla^2(\ell)(\beta) u + \lambda^T \hat{\rho} u$$

for some $\lambda \in \mathbb{R}^{d+1}$ and is characterized by

$$\begin{cases} \nabla^2(\ell)(\beta)u + \nabla\ell(\beta) + \hat{\rho}^T\lambda = 0 \\ \qquad\qquad\qquad\quad \hat{\rho}u = 0 \end{cases}$$

This shows that the Newton-Raphson iterations can be implemented as

$$\beta_{n+1} = \beta_n - \epsilon_{n+1}u_{n+1} \tag{8.6}$$

with

$$\begin{pmatrix} u_{n+1} \\ \lambda \end{pmatrix} = \begin{pmatrix} \nabla^2(\ell)(\beta_n) & \hat{\rho}^T \\ \hat{\rho} & 0 \end{pmatrix}^{-1} \begin{pmatrix} \nabla\ell(\beta_n) \\ 0 \end{pmatrix}. \tag{8.7}$$

We summarize this discussion in the following algorithm.

---

**Algorithm 8.1 (Logistic regression with Newton's gradient ascent)**

(1) Input: (i) training data $(x_1, y_1, \ldots, x_N, y_N)$ with $x_i \in \mathbb{R}^d$ and $y_i \in \mathcal{R}_Y$; (ii) coefficients $\rho_g, g \in \mathcal{R}_Y$ with non-zero sum and target value $c \in \mathbb{R}$; (iii) algorithm step $\epsilon$ small enough.

(2) Initialize the algorithm with $\beta_0$ such that $\sum_g \rho_g \beta_0(g) = c$.

(3) At iteration $n$, compute $\nabla\ell(\beta_n)$ and $\nabla^2(\ell)(\beta_n)$ as provided by proposition 8.1.

(4) Update $\beta_n$ using (8.6) and (8.7), with $\epsilon_{n+1} = \epsilon$. Alternatively, optimize $\epsilon_{n+1}$ using a line search.

(5) Stop the procedure if the change in the parameter is below a small tolerance level. Otherwise, return to step 2.

---

### 8.1.4 Penalized Logistic Regression

Logistic regression can be combined with a penalty term, e.g., maximizing

$$\ell_2(\beta) = \ell(\beta) - \lambda \sum_{i=1}^{d} |b^{(i)}|^2 \tag{8.8}$$

or

$$\ell_1(\beta) = \ell(\beta) - \lambda \sum_{i=1}^{d} |b^{(i)}| \tag{8.9}$$

where $b^{(i)}$ is the $q$-dimensional vector formed with the $i$th coefficients of $b(g)$ for $g \in \mathcal{R}_Y$. Similarly to penalized regression, one generally normalizes the $x$ variables to have unit standard deviation before applying the method.

**Maximization with the $\ell_2$ norm**   The problem in (8.8) relates to ridge regression and can be solved using a Newton-Raphson method (Algorithm 8.1) with minor changes. More precisely, letting

$$\Delta = \begin{pmatrix} 0 & 0 \\ 0 & \mathrm{Id}_{\mathbb{R}^d} \end{pmatrix}$$

we have, considering $\beta$ as a $d+1$ by $q$ matrix,

$$\ell_2(\beta) = \ell(\beta) - \lambda \mathrm{trace}(\beta^T \Delta \beta)$$

and

$$d\ell_2(\beta)u = d\ell(\beta)u - 2\lambda \mathrm{trace}(\beta^T \Delta u),$$

$$d^2\ell_2(\beta)(u,u') = d\ell(\beta)(u,u') - 2\lambda \mathrm{trace}(u^T \Delta u').$$

In addition, when $\lambda > 0$, the problem is over-parametrized only up to the addition of a constant to $(g \mapsto a_0(g))$, so that one only needs a single constraint $\sum_g \rho_g a_0(g) = 0$ and the Lagrange coefficient in (8.7) is one dimensional.

**Maximization with the $\ell_1$ norm**   The maximization in (8.9) can be run using proximal gradient ascent (section 3.5.5), writing the objective function in the form

$$\ell_1(a_0, b) = \ell(a_0, b) - \lambda \gamma(a_0, b)$$

with

$$\gamma(a_0, b) = \sum_{i=1}^d \sqrt{\sum_{g \in \mathcal{R}_Y} b^{(i)}(g)^2}.$$

Here, $\ell$ is concave and $\gamma$ is convex and the proximal gradient iterations are

$$\beta_{n+1} = \mathrm{prox}_{\epsilon \lambda \gamma}(\beta_n + \epsilon \overline{\nabla}\ell(\beta_n)) \tag{8.10}$$

where $\overline{\nabla}\ell$ is the gradient of $\ell$ projected on the set of functions $u : \mathcal{R}_Y \to \mathbb{R}^{d+1}$ satisfying

$$\sum_{g \in \mathcal{R}_Y} \rho_g u^{(0)}(g) = 0$$

where $u^{(0)}(g)$ is the first coordinate of $u(g)$. This projection can be computed by subtracting

$$\frac{\sum_{g' \in \mathcal{R}_Y} \rho(g') \partial_{a_0(g')} \ell}{\sum_{g' \in \mathcal{R}_Y} \rho(g')^2} \rho(g)$$

to $\partial_{a_0(g)}\ell$. This algorithm will converge for small enough $\epsilon$.

We already know the gradient of $\ell$, so it only remains to determine the proximal operator of $\gamma$ to make (8.10) explicit. Let us denote the coordinates of a function $u : \mathcal{R}_Y \to \mathbb{R}^{d+1}$ as $u^{(i)}(g)$ for $i = 0, \ldots, d$ and $g \in \mathcal{R}_Y$.

$$\text{prox}_{\epsilon\lambda\gamma}(u) = \underset{\tilde{u}}{\text{argmin}} \left( \gamma(\tilde{u}) + \frac{1}{2\lambda\epsilon} \sum_{i=0}^{d} \sum_{g \in \mathcal{R}_Y} (u^{(i)}(g) - \tilde{u}^{(i)}(g))^2 \right).$$

Since $\gamma$ does not depend on $\tilde{u}^{(0)}(\cdot)$, the optimal $\tilde{u}^{(0)}(\dot{)}$ is $\tilde{u}^{(0)}(\cdot) = u^{(0)}(\cdot)$. One can optimize separately each group $(\tilde{u}^{(i)}(g), g \in \mathcal{R}_Y)$, which must minimize

$$\sqrt{\sum_{g \in \mathcal{R}_Y} \tilde{u}^{(i)}(g)^2} + \frac{1}{2\lambda\epsilon} \sum_{g \in \mathcal{R}_Y} (u^{(i)}(g) - \tilde{u}^{(i)}(g))^2.$$

The function $t \mapsto \sqrt{t}$ being differentiable everywhere except at 0, we first search for a solution for which at least one $\tilde{u}^{(i)}(g)$ does not vanish. If such a solution exists, it must satisfy, for all $g \in \mathcal{R}_Y$

$$\frac{\tilde{u}^{(i)}(g)}{\sqrt{\sum_{g' \in \mathcal{R}_Y} \tilde{u}^{(i)}(g')^2}} + \frac{1}{\lambda\epsilon}(\tilde{u}^{(i)}(g) - u^{(i)}(g)) = 0$$

Letting $|\tilde{u}^{(i)}(\cdot)| = \sqrt{\sum_{g \in \mathcal{R}_Y} \tilde{u}^{(i)}(g)^2}$ we get

$$\tilde{u}^{(i)}(\cdot)(|\tilde{u}^{(i)}(\cdot)| + \lambda\epsilon) = u^{(i)}(\cdot)|\tilde{u}^{(i)}(\cdot)|$$

Taking the norm on both sides and dividing by $|\tilde{u}(i, \cdot)|$ (which is assumed not to vanish) yields

$$|\tilde{u}^{(i)}(\cdot)| + \lambda\epsilon = |u^{(i)}(\cdot)|,$$

which has a positive solution only if $|u^{(i)}(\cdot)| > \lambda\epsilon$, and gives in that case

$$\tilde{u}^{(i)}(\cdot) = \frac{|u^{(i)}(\cdot)| - \epsilon\lambda}{|u^{(i)}(\cdot)|} u^{(i)}(\cdot)$$

If $|u^{(i)}(\cdot)| \leq \lambda\epsilon$, then we must take $\tilde{u}^{(i)}(\cdot) = 0$. We have therefore obtained:

$$\text{prox}_{\epsilon\lambda g}(a) = \tilde{u}$$

with $\tilde{u}^{(0)}(\cdot) = u^{(0)}(\cdot)$ and

$$\tilde{u}^{(i)}(\cdot) = \max\left( \frac{|u^{(i)}(\cdot)| - \epsilon\lambda}{|u^{(i)}(\cdot)|}, 0 \right) u^{(i)}(\cdot)$$

for $i \geq 1$. We summarize this discussion in the next algorithm.

**Algorithm 8.2 (Logistic lasso)**
(1) Input: (i) training data $(x_1, y_1, \ldots, x_N, y_N)$ with $x_i \in \mathbb{R}^d$ and $y_i \in \mathcal{R}_Y$; (ii) coefficients $\rho_g, g \in \mathcal{R}_Y$ with non-zero sum and target value $c \in \mathbb{R}$; (iii) algorithm step $\epsilon$; (iv) penalty coefficient $\lambda$.

(2) Initialize the algorithm with $\beta_0 = (a_{00}, b_0)$ such that $\sum_g \rho_g a_{00}(g) = c$.

(3) At iteration $n$, compute $u = \beta_n + \epsilon \overline{\nabla} \ell(\beta_n)$, with $\beta_n = (a_{0,n}, b_n)$.

(4) Let $a_{n+1,0}(\cdot) = u^{(0)}(\cdot)$ and for $i \geq 1$,

$$b_{n+1}^{(i)}(\cdot) = \max\left(\frac{|u^{(i)}(\cdot)| - \epsilon\lambda}{|u^{(i)}(\cdot)|}, 0\right) u^{(i)}(\cdot)$$

(5) Stop the procedure if the change in the parameter is below a small tolerance level. Otherwise, return to step 2.

---

### 8.1.5   Kernel logistic regression

Let $h : \mathcal{R}_X \to H$ be a feature function with values in a Hilbert space $H$ with $K(x, x') = \langle h(x), h(x') \rangle_H$. The kernel version of logistic regression uses the model:

$$\log p_{a_0, b}(g \mid x) = a_0(g) + \langle h(x), b(g) \rangle_H - \log \sum_{\tilde{g} \in \mathcal{R}_Y} \exp(a_0(\tilde{g}) + \langle h(x), b(\tilde{g}) \rangle_H)$$

with $b(g) \in H$ for $g \in \mathcal{R}_Y$.

Using the usual kernel argument, one sees that, when maximizing the log-likelihood, there is no loss of generality is assuming that each $b(g)$ belongs to $V = \text{span}(h(x_1), \ldots, h(x_N))$. Taking

$$b(g) = \sum_{k=1}^{N} \alpha_k(g) h(x_k),$$

we have

$$\log p_\alpha(g \mid x) = a_0(g) + \sum_{k=1}^{N} \alpha_k(g) K(x, x_k) - \log\left(\sum_{\tilde{g} \in \mathcal{R}_Y} \exp(a_0(\tilde{g}) + \sum_{k=1}^{N} \alpha_k(\tilde{g}) K(x, x_k))\right).$$

To avoid overfitting, one must include a penalty term in the likelihood, and (in order to take advantage of the kernel), one can take this term proportional to $\sum_g \|b(g)\|_H^2$. The complete learning procedure then requires to maximize the concave penalized likelihood

$$\ell(\alpha) = \sum_{k=1}^{N} \log p_\alpha(y_k \mid x_k) - \lambda \sum_{g \in \mathcal{R}_Y} \sum_{k,l=1}^{N} \alpha_k(g) \alpha_l(g) K(x_k, x_l).$$

The computation of the first and second derivatives of this function is similar to that for the original version, and we skip the details.

## 8.2 Linear Discriminant analysis

### 8.2.1 Generative model in classification and LDA

**Generative model**  In classification, the class variable $Y$ generally has a causal role upon which the variable $X$ is produced. Prediction can therefore be seen as an inverse problem where the cause is deduced from the result. In terms of generative modeling, one should therefore model the distribution of $Y$, followed by the the conditional distribution of $X$ given $Y$.

Taking $\mathcal{R}_X = \mathbb{R}^d$, denote by $f_g$ the conditional p.d.f. of $X$ given $Y = g$ and let $\pi_g = P(Y = g)$. The Bayes estimator for the 0–1 loss maximizes the posterior probability

$$\mathbb{P}(Y = g \mid X = x) = \frac{\pi_g f_g(x)}{\sum_{g' \in \mathcal{R}_Y} \pi_{g'} f_{g'}(x)}.$$

Since the denominator does not depend on $g$ the Bayes estimator equivalently maximizes (taking logarithms)

$$\log f_g(x) + \log \pi_g.$$

One generally speaks of a linear classification method when the prediction is based on the maximization in $g$ of a function $U(g,x)$ where $U$ is affine in $x$. In this sense, logistic regression is linear, and kernel logistic regression is linear in feature space. For the generative approach, this occurs when one uses the following model, which provides the generative form of linear discriminant analysis (LDA). Assume that the distributions $f_g$ are all Gaussian with mean $m_g$ and *common variance S*, so that

$$f_g(x) = \frac{1}{\sqrt{(2\pi)^d \det \Sigma}} e^{-\frac{1}{2}(x-m_g)^T S^{-1}(x-m_g)}. \tag{8.11}$$

In this case, the optimal predictor must maximize (in $g$)

$$-\frac{1}{2}(x - m_g)^T S^{-1}(x - m_g) + \log \pi_g.$$

Introduce $m = \mathbb{E}(X) = \sum_{g \in \mathcal{R}_Y} \pi_g m_g$. Then the optimal classifier must maximize

$$-\frac{1}{2}(x - m)^T S^{-1}(x - m) + (x - m)^T S^{-1}(m_g - m) - \frac{1}{2}(m_g - m)^T S^{-1}(m_g - m) + \log \pi_g.$$

Since the first term does not depend on $g$, it is equivalent to maximize

$$(x - m)^T S^{-1}(m_g - m) - \frac{1}{2}(m_g - m)^T S^{-1}(m_g - m) + \log \pi_g \tag{8.12}$$

with respect to the class $g$, which provides an affine function of $x$.

**Training**   Training for LDA simply consists in estimating the class means and common variance in (8.11) from data. We introduce some notation for this purpose (this notation will be reused through the rest of this chapter).

Recall that $N_g, g \in \mathcal{R}_Y$ denotes the number of samples with class $g$ in the training set $T = (x_1, y_1, \ldots, x_N, y_N)$. We let $c_g = N_g/N$ and $C$ be the diagonal matrix with diagonal coefficients $c_{g_1}, \ldots, c_{g_q}$. We also let $\zeta \in \mathbb{R}^q$ denote the vector with the same coordinates. For $g \in \mathcal{R}_Y$, $\mu_g$ denotes the class average

$$\mu_g = \frac{1}{N_g} \sum_{k=1}^{n} x_k \mathbf{1}_{y_k = g}$$

and $\mu$ the global average

$$\mu = \frac{1}{N} \sum_{k=1}^{N} x_k = \sum_{g \in \mathcal{R}_Y} c_g \mu_g.$$

Let $\Sigma_g$ denote the sample covariance matrix in class $g$, defined by

$$\Sigma_g = \frac{1}{N_g} \sum_{k=1}^{N} (x_k - \mu_g)(x_k - \mu_g)^T \mathbf{1}_{y_k = g},$$

and $\Sigma_w$ the pooled class covariance (also called *within-class* covariance) defined by

$$\Sigma_w = \frac{1}{N} \sum_{k=1}^{N} (x_k - \mu_{y_k})(x_k - \mu_{y_k})^T = \sum_{g \in \mathcal{R}_Y} c_g \Sigma_g.$$

Let, in addition, $\Sigma_b$ denotes the "between-class" covariance matrix, given by

$$\Sigma_b = \sum_{g \in \mathcal{R}_Y} c_g (\mu_g - \mu)(\mu_g - \mu)^T$$

The global covariance matrix, given by,

$$\Sigma_{XX} = \frac{1}{N} \sum_{k=1}^{N} (x_k - \mu)(x_k - \mu)^T$$

satisfies $\Sigma_{XX} = \Sigma_w + \Sigma_b$. This identity is proved by noting that, for any $g \in \mathcal{R}_Y$,

$$\frac{1}{N_g} \sum_{k=1}^{N} (x_k - \mu)(x_k - \mu)^T \mathbf{1}_{y_k = g} = \Sigma_g + (\mu_g - \mu)(\mu_g - \mu)^T.$$

We will finally denote by $M$ the matrix

$$M = \begin{pmatrix} (\mu_{g_1} - \mu)^T \\ \vdots \\ (\mu_{g_q} - \mu)^T \end{pmatrix}.$$

Note that $\Sigma_b = M^T C M$.

Given this notation, one can in particular take $\hat{m}_g = \mu_g$, $m = \mu$ and $\hat{S} = \Sigma_w$ in (8.12). The class probabilities, $\pi_g$, can be deduced from the normalized frequencies of $y_1, \ldots, y_N$. However, in many applications, one prefers to simply fix $\pi_g = 1/q$, in order to balance the importance of each class.

**Remark 8.5** If one relaxes the assumption of common class variances, one needs to use $\Sigma_g$ in place of $\Sigma_w$ for class $g$. The decision boundaries are not linear in this case, but provided by quadratic equations (and the resulting method if often called quadratic discriminant analysis, or QDA). QDA requires the estimation of $qd(d + 3)/2$ coefficients, which may be overly ambitious when the sample size is not large compared to the dimension, in which case QDA is prone to overfitting. (Even LDA, which involves $qd + d(d+1)/2$ parameters, may be unrealistic in some cases.) We also note a variant of QDA that uses class covariance matrices given by

$$\tilde{\Sigma}_g = \alpha \Sigma_w + (1 - \alpha)\Sigma_g.$$

### 8.2.2 Dimension reduction

One of the interests of LDA is that it can be combined with a rank reduction procedure. LDA with $q$ classes can always be seen as a $(q-1)$-dimensional problem after suitable projection on a data-dependent affine space. Recall that the classification rule after training requires to maximize w.r.t. $g \in \mathcal{R}_Y$ the function

$$(x - \mu)^T \Sigma_w^{-1}(\mu_g - \mu) - \frac{1}{2}(\mu_g - \mu)^T \Sigma_w^{-1}(\mu_g - \mu) + \log \pi_g.$$

Define the "spherized" data [2] by $\tilde{x}_k = \Sigma_w^{-1/2}(x_k - \mu)$, where $\Sigma_w^{1/2}$ is the positive symmetric square root of $\Sigma_w$. Also let $\tilde{\mu}_g = \Sigma_w^{-1/2}(\mu_g - \mu)$.

With this notation, the predictor chooses the class $g$ that maximizes

$$\tilde{x}^T \tilde{\mu}_g - \frac{1}{2}|\tilde{\mu}_g|^2 + \log \pi_g$$

with $\tilde{x} = \Sigma_w^{-1/2}(x - \bar{\mu})$.

---

[2]In this section only, the notation $\tilde{x}$ does not refer to $(1, x^T)^T$.

Now, let $V = \mathrm{span}\{\tilde{\mu}_g, g \in \mathcal{R}_Y\}$. Since $\sum_g c_g \tilde{\mu}_g = 0$, this space is at most $(q-1)$-dimensional. Let $P_V$ denote the orthogonal projection on $V$. We have $\tilde{x}^T z = (P_V \tilde{x})^T z$ for any $z \in V$ and $\tilde{x} \in \mathbb{R}^d$.

The classification rule can then be replaced by maximizing

$$(P_V \tilde{x})^T \tilde{\mu}_g - \frac{1}{2}|\tilde{\mu}_g|^2 + \log \pi_g$$

with $\tilde{x} = \Sigma_w^{-1/2}(x - \bar{\mu})$.

Recall that $M = \begin{pmatrix} (\mu_{g_1} - \mu)^T \\ \vdots \\ (\mu_{g_q} - \mu)^T \end{pmatrix}$ and let $\widetilde{M} = \begin{pmatrix} \tilde{\mu}_{g_1}^T \\ \vdots \\ \tilde{\mu}_{g_q}^T \end{pmatrix}$. The dimension, denoted $r$, of $V$ is equal to the rank of $\widetilde{M}$. Let $(\tilde{e}_1, \ldots, \tilde{e}_r)$ be an orthonormal basis of $V$. One has

$$P_V \tilde{x} = \sum_{j=1}^{r} (\tilde{x}^T \tilde{e}_j)\tilde{e}_j.$$

Given an input $x$, one must therefore compute the "scores" $\gamma_j(x) = \tilde{x}^T \tilde{e}_j$ and maximize

$$\sum_{j=1}^{r} \gamma_j(x)\gamma_j(\mu_g) - \frac{1}{2}\sum_{j=1}^{r} \gamma_j(\mu_g)^2 + \log \pi_g.$$

The following proposition is key to the practical implementation of LDA with dimension reduction.

**Proposition 8.6** *An orthonormal basis of $V = \mathrm{span}(\tilde{\mu}_g, g \in CG)$ is provided by the the first $r$ eigenvectors of $\widetilde{M}^T C\widetilde{M}$ associated with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_r > 0$ (all other eigenvalues being zero).*

Proof  Indeed, if $\tilde{x}$ is perpendicular to $V$, we have

$$\widetilde{M}^T C\widetilde{M}\tilde{x} = \sum_{g \in \mathcal{R}_Y} c_g(\tilde{\mu}_g^T \tilde{x})\tilde{\mu}_g = 0$$

so that $V^\perp \subset \mathrm{Null}(\widetilde{M}^T C\widetilde{M})$, and both spaces coincide because they have the same dimension $(d - r)$. This shows that $V = \mathrm{Null}(\widetilde{M}^T C\widetilde{M})^\perp = \mathrm{Range}(\widetilde{M}^T C\widetilde{M})$. Since $\widetilde{M}^T C\widetilde{M}$ is symmetric, $\mathrm{Null}(\widetilde{M}^T C\widetilde{M})^\perp$ is generated by eigenvectors with non-zero eigenvalues.                                                                      ∎

Returning to the original variables, we have $\widetilde{M} = M\Sigma_w^{-1/2}$ and $M^T CM = \Sigma_b$, the between class covariance matrix. This implies that $\widetilde{M}^T C\widetilde{M} = \Sigma_w^{-1/2}\Sigma_b\Sigma_w^{-1/2}$ and each

Figure 8.1: Left: Original (training) data with three classes. Right: LDA scores, where the $x$ axis provides $\gamma_1$ and the $y$ axis $\gamma_2$.

eigenvector $\tilde{e}_j$ therefore satisfies

$$\Sigma_b \Sigma_w^{-1/2} \tilde{e}_j = \lambda_j \Sigma_w^{1/2} \tilde{e}_j = \lambda_j \Sigma_w (\Sigma_w^{-1/2} \tilde{e}_j).$$

Therefore, letting $e_j = \Sigma_w^{-1/2} \tilde{e}_j$, $(e_1, \ldots, e_r)$ are the solutions of the generalized eigenvalue problem $\Sigma_b e = \lambda \Sigma_w e$ that are associated with non-zero eigenvalues (they are however normalized so that $e_j^T \Sigma_w e_j = 1$). Moreover, the scores are given by

$$\gamma_j(x) = \tilde{x}^T \tilde{e}_j = (x - \mu)^T \Sigma_w^{-1/2} \tilde{e}_j = (x - \mu)^T e_j$$

and can therefore be computed directly from the original data and the vectors $e_1, \ldots, e_r$. An example of training data and its representation in the LDA space (associated with the scores) in provided in fig. 8.1.

We can now describe the LDA learning algorithm with dimension reduction.

---

**Algorithm 8.3 (LDA with dimension reduction)**

1. Compute $\mu_g, g \in \mathcal{R}_Y$, $\Sigma_w$ and $\Sigma_b$ from training data.

2. Estimate (if needed) $\pi_g, g \in \mathcal{R}_Y$

3. Solve the generalized eigenvalue problem $\Sigma_b e = \lambda \Sigma_w e$. Let $e_1, \ldots, e_r$ be the eigenvectors associated with non-zero eigenvalues, normalized so that $e_j^T \Sigma_w e_j = 1$.

4. Choose a reduced dimension $r_0 \le r$.

5. Precompute mean scores $\gamma_j(\mu_g) = (\mu_g - \mu)^T e_j$, $g \in \mathcal{R}_Y, j = 1, \ldots, r_0$.

6. To classify a new example $x$, compute $\gamma_j(x) = (x - \mu)^T e_j$ and choose the class that maximizes

$$\sum_{j=1}^{r_0} \gamma_j(x) \gamma_j(\mu_g) - \frac{1}{2} \sum_{j=1}^{r_0} \gamma_j(\mu_g)^2 + \log \pi_g.$$

---

### 8.2.3   Fisher's LDA

This characterization leads to the discriminative interpretation of LDA, also called Fisher's LDA. Indeed, the generalized eigenvalue problem $\Sigma_b e = \lambda \Sigma_w e$ is directly related to the maximization of the ratio $e^T \Sigma_b e$ subject to $e^T \Sigma_w e = 1$, which provides directions that have a large between-class variance for within class variance equal to 1. More precisely, $e_1$ is the direction that achieves the maximum; $e_2$ is the second best direction, constrained to being perpendicular to $e_1$, and so on until $e_r$ which is the optimal constrained to be perpendicular to $(e_1, \ldots, e_{r-1})$. We are therefore looking for directions that have the largest ratio of between-class variance to within-class variance.

### 8.2.4   Kernel LDA

**Mean and covariance in feature space**    We assume the usual construction where $h : \mathcal{R}_X \to H$ is a feature function, $H$ a Hilbert space with kernel $K(x, x') = \langle h(x), h(x') \rangle_H$. (The assumption that $H$ is a complete space is here required for the expectations below to be meaningful.)

We now discuss the kernel version of LDA by plugging the feature space representation directly in the classification rule. So, consider $h : \mathcal{R} \to H$. Let $X : \Omega \to \mathcal{R}$ be a random variable such that $\mathbb{E}(\|h(X)\|_H^2) < \infty$. Then, its mean feature $m = \mathbb{E}(h(X))$ is well defined as an element of $H$, and so are the class averages, $m_g = \mathbb{E}(h(X) \mid Y = g)$.

In this possibly infinite-dimensional setting, the covariance "matrix" is defined as a linear operator $S : H \to H$ such that, for all $\xi, \eta \in H$:

$$\langle \xi, S\eta \rangle_H = \mathbb{E}\left( \langle h(X) - m, \xi \rangle_H \langle h(X) - m, \eta \rangle_H \right), \tag{8.13}$$

which is equivalent to defining

$$S\eta = E(\langle h(X) - m, \eta \rangle_H (h(X) - m))$$

for $\eta \in H$. This definition generalizes the identity for a random variable $U : \Omega \to \mathbb{R}^d$:

$$S_U w = \mathbb{E}((U - E(U))(U - E(U))^T)w = \mathbb{E}(((U - E(U))^T w)(U - E(U)))$$

One can similarly define the covariance matrix in class $g$, $S_g$, by conditioning the right-hand side in (8.13) by $Y = g$ and replacing $m$ by $m_g$.

**LDA in feature space**

Following the LDA model, we assume that the operators $S_g$ are all equal to a fixed operator, the within-class covariance operator denoted $S$.

Assuming that $S$ is invertible, one can generalize the LDA classification rule to data represented in feature space by classifying a new input $x$ in class $g$ when

$$\langle h(x) - m, S^{-1}(m_g - m)\rangle_H - \frac{1}{2}\langle m_g - m, S^{-1}(m_g - m)\rangle_H + \log \pi_g \qquad (8.14)$$

is maximal over all classes. Notice that this is a transcription of the finite-dimensional Bayes rule, but cannot be derived from a generative model, because the assumption that $h(X)$ is Gaussian is not valid in general. (It would require that $h$ takes values in a $d$-dimensional linear space, which would eliminate all interesting kernel representations.)

Let, as before, $T = (x_1, y_1, \ldots, x_N, y_N)$ be the training set, $N_g$ denote the number of examples in class $g$ and $c_g = N_g/N$. When $h$ is known (which, we recall, is not a practical assumption, but we will fix this later), one can estimate the class averages from training data by

$$\mu_g = \frac{1}{N_g}\sum_{k=1}^{N} h(x_k)\mathbf{1}_{y_k=g}$$

and the within-class covariance operator by

$$\langle \xi, \Sigma_w \eta\rangle_H = \frac{1}{N}\sum_{k=1}^{N} \langle h(x_k) - \mu_{y_k}, \xi\rangle_H \langle h(x_k) - \mu_{y_k}, \eta\rangle_H.$$

Unfortunately, the resulting variance estimator cannot be directly used in (8.14), because it is not invertible if $\dim(H) > N$. Indeed, one has $\Sigma_w \eta = 0$ as soon as $\eta$ is perpendicular to $V \stackrel{\Delta}{=} \mathrm{span}(h(x_1), \ldots, h(x_N))$.

One way to address the degeneracy of the estimated covariance operator is to add to $\Sigma_w$ a small multiple of the identity, say $\rho\mathrm{Id}_H$,[3] and let the classification rule maximize in $g$:

$$\langle h(x) - \mu, (\Sigma_w + \rho\mathrm{Id}_H)^{-1}(\mu_g - \mu)\rangle_H - \frac{1}{2}\langle \mu_g - \mu, (\Sigma_w + \rho\mathrm{Id}_H)^{-1}(\mu_g - \mu)\rangle_H + \log \pi_g.$$
$$(8.15)$$

where $\mu$ is the average of $h(x_1), \ldots, h(x_N)$. Taking this option, we still need to make this expression computable and remove the dependency in the feature function $h$.

**Reduction**   We have $\mu_g \in V$ for all $g \in \mathcal{R}_Y$ and, since

$$\Sigma_w \eta = \frac{1}{N}\sum_{k=1}^{N} \langle h(x_k) - \mu_{y_k}, \eta\rangle_H (h(x_k) - \mu_{y_k}),$$

---

[3]The operator $A + \rho\mathrm{Id}_H$ is invertible as soon as $A$ is symmetric positive semi-definite.

this operator maps $H$ to $V$, which implies that $\Sigma_w + \rho \mathrm{Id}_H$ maps $V$ into itself. Moreover, this mapping is onto: If $v \in V$ and $u = (\Sigma_w + \rho \mathrm{Id}_H)^{-1} v$, then, $u \in V$. Indeed, for any $z \perp V$, we have $\langle z, \Sigma_w u + \rho u \rangle_H = \langle z, v \rangle_H$. We have $\langle z, \Sigma_w u \rangle_H = 0$ (because $\Sigma_w$ maps $H$ to $V$) and $\langle z, v \rangle_H = 0$ (because $v \in V$), so that we can conclude that $\langle z, u \rangle_H = 0$. Since this is true for all $z \perp V$, this requires that $u \in V$.[4]

We now express the classification rule in (8.15) as a function of the kernel associated with the feature-space representation. Denote, for any vector $u \in \mathbb{R}^N$,

$$\xi(u) = \sum_{k=1}^{N} u^{(k)} h(x_k),$$

therefore defining a mapping $\xi$ from $\mathbb{R}^N$ onto $V$. Letting as usual $\mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$ be the matrix formed by pairwise evaluations of $K$ on training inputs, we have the identity

$$\langle \xi(u), \xi(u') \rangle_H = u^T \mathcal{K} u'.$$

for all $u, u' \in \mathbb{R}^N$. For simplicity, we will assume in the rest of the discussion that $\mathcal{K}$ is invertible.

We have $\mu_g = \xi(\mathbb{1}_g / N_g)$, where $\mathbb{1}_g \in \mathbb{R}^N$ is the vector with $k$th coordinate equal to 1 if $y_k = g$ and 0 otherwise. Also $\mu = \xi(\mathbb{1}/N)$ (recall that $\mathbb{1}$ is the vector with all coordinates equal to 1).

For $u \in \mathbb{R}^N$, we want to characterize $v \in \mathbb{R}^N$ such that $\Sigma_w \xi(u) = \xi(v)$. Let $\delta_k$ denote the vector with 1 at the $k$th entry and 0 otherwise. We have

$$\Sigma_w \xi(u) = \frac{1}{N} \sum_{k=1}^{N} \langle \xi(u), h(x_k) - \mu_{y_k} \rangle_H (h(x_k) - \mu_{y_k})$$

$$= \frac{1}{N} \sum_{k=1}^{N} \langle \xi(u), \xi(\delta_k - \mathbb{1}_{y_k}/N_{y_k}) \rangle_H \xi(\delta_k - \mathbb{1}_{y_k}/N_{y_k})$$

$$= \frac{1}{N} \sum_{k=1}^{N} ((\delta_k - \mathbb{1}_{y_k}/N_{y_k})^T \mathcal{K} u) \, \xi(\delta_k - \mathbb{1}_{y_k}/N_{y_k})$$

$$= \xi \left( \frac{1}{N} \sum_{k=1}^{N} ((\delta_k - \mathbb{1}_{y_k}/N_{y_k})^T \mathcal{K} u)(\delta_k - \mathbb{1}_{y_k}/N_{y_k}) \right)$$

so that $\Sigma_w \xi(u) = \xi(P \mathcal{K} u)$ with

$$P = \frac{1}{N} \sum_{k=1}^{N} (\delta_k - \mathbb{1}_{y_k}/N_{y_k})(\delta_k - \mathbb{1}_{y_k}/N_{y_k})^T$$

---

[4] One has $(V^\perp)^\perp = V$ for finite-dimensional—or more generally closed—subspaces of $H$

Note that one has

- $$\sum_{k=1}^{N} \delta_k \delta_k^T = \mathrm{Id}_{\mathbb{R}^N},$$

- $$\sum_{k=1}^{N} \left(\frac{\mathbb{1}_{y_k}}{N_{y_k}}\right) \delta_k^T = \sum_{g \in \mathcal{R}_Y} \frac{\mathbb{1}_g}{N_g} \sum_{k: y_k = g} \delta_k^T = \sum_{g \in \mathcal{R}_Y} \frac{\mathbb{1}_g \mathbb{1}_g^T}{N_g} = \sum_{k=1}^{N} \delta_k \left(\frac{\mathbb{1}_{y_k}}{N_{y_k}}\right)^T,$$

- $$\sum_{k=1}^{N} \left(\frac{\mathbb{1}_{y_k}}{N_{y_k}}\right) \left(\frac{\mathbb{1}_{y_k}}{N_{y_k}}\right)^T = \sum_{g \in \mathcal{R}_Y} \frac{\mathbb{1}_g \mathbb{1}_g^T}{N_g}.$$

This shows that $P$ can be expressed as

$$P = \frac{1}{N}\left(\mathrm{Id}_{\mathbb{R}^N} - \sum_{g \in \mathcal{R}_Y} \mathbb{1}_g \mathbb{1}_g^T / N_g\right).$$

We have therefore proved that

- $(\Sigma_w + \rho \mathrm{Id}_H)\xi(u) = \xi\left((P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})u\right)$
- $(\Sigma_w + \rho \mathrm{Id}_H)^{-1}\xi(\tilde{u}) = \xi\left((P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1}\tilde{u}\right).$

Recall that the feature-space LDA classification rule maximizes

$$\langle h(x) - \mu, (\Sigma_w + \rho \mathrm{Id}_H)^{-1}(\mu_g - \mu)\rangle_H - \frac{1}{2}\langle \mu_g - \mu, (\Sigma_w + \rho \mathrm{Id}_H)^{-1}(\mu_g - \mu)\rangle_H + \log \pi_g.$$

All terms belong to $V$, except $h(x)$, but this term can be replaced by its orthogonal projection on $V$ without changing the result. This projection can be made explicit in terms of the representation $\xi$ as follows. For $x \in \mathcal{R}$, let $\xi(\psi(x))$ denote the orthogonal projection of $h(x)$ on $V$ (this defines the function $\psi$). If $v(x)$ denotes the vector with coordinates $K(x, x_k)$, $k = 1, \ldots, N$, then $\psi(x) = \mathcal{K}^{-1} v(x)$, as can be obtained by identifying the inner products $\langle h(x), h(x_k)\rangle_H$ and $\langle \xi(\psi(x)), h(x_k)\rangle_H$.

We are now ready to rewrite the kernel LDA classification rule in terms of quantities that only involve $K$. We have

$$\langle h(x) - \mu, (\Sigma_w + \rho \mathrm{Id}_H)^{-1}(\mu_g - \mu)\rangle_H$$
$$= \langle \xi(\psi(x) - \mathbb{1}/N), \xi((P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1}(\mathbb{1}_g/N_g - \mathbb{1}/N))\rangle_H$$
$$= (\psi(x) - \mathbb{1}/N)^T \mathcal{K}(P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1}(\mathbb{1}_g/N_g - \mathbb{1}/N)$$

Given this, the classification rule must maximize

$$(\psi(x) - \mathbb{1}/N)^T \mathcal{K}(P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1}(\mathbb{1}_g/N_g - \mathbb{1}/N)$$
$$- \frac{1}{2}(\mathbb{1}_g/N_g - \mathbb{1}/N)^T \mathcal{K}(P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1}(\mathbb{1}_g/N_g - \mathbb{1}/N) + \log \pi_g. \quad (8.16)$$

**Dimension reduction**   Note that $\mathcal{K}(P\mathcal{K} + \rho \mathrm{Id}_{\mathbb{R}^N})^{-1} = \mathcal{K}(\mathcal{K}P\mathcal{K} + \rho\mathcal{K})^{-1}\mathcal{K}$ is a symmetric matrix. So, the expression in (8.16) can be written as

$$(v(x) - v)^T R^{-1}(v_g - v) - \frac{1}{2}(v_g - v)^T R^{-1}(v_g - v) + \log \pi_g.$$

with $R = \mathcal{K}P\mathcal{K} + \rho\mathcal{K}$, $v_g = \mathcal{K}\mathbb{1}_g/N_g$ and $\bar{v} = \mathcal{K}\mathbb{1}/N$. Clearly, if $v_1,\dots,v_N$ are the column vectors $\mathcal{K}$, we have

$$v_g = \frac{1}{N_g}\sum_{k=1}^{N} v_k \mathbf{1}_{y_k=g}, \quad \bar{v} = \frac{1}{N}\sum_{k=1}^{N} v_k.$$

We therefore retrieve an expression similar to finite-dimensional LDA, provided that one replaces $x$ by $v(x)$, $x_k$ by $v_k$ and $\Sigma_w$ by $R$. Letting

$$Q = \frac{1}{N}\sum_{g \in \mathcal{R}_Y} N_g (v_g - \bar{v})(v_g - \bar{v})^T$$

be the between-class covariance matrix, the discriminant directions are therefore solutions of the generalized eigenvalue problem

$$Qf_j = \lambda_j Rf_j$$

with $f_j^T R f_j = 1$ with $R = (\mathcal{K}P\mathcal{K} + \rho\mathcal{K})$. Note that

$$\mathcal{K}P\mathcal{K} = \frac{1}{N}\sum_{k=1}^{N}(v_k - \bar{v}_{y_k})(v_k - \bar{v}_{y_k})^T$$

is the within-class covariance matrix for the training data $(v_1, y_1, \dots, v_N, y_N)$.

The following summarizes the kernel LDA classification algorithm.

---

**Algorithm 8.4 (Kernel LDA)**

(1) Select a positive kernel $K$ and a coefficient $\rho > 0$.

(2) Given $T = (x_1, y_1, \dots, x_N, y_N)$ and a kernel $K$, compute the kernel matrix $\mathcal{K} = \mathcal{K}(x_1, \dots, x_N)$ and the matrix $R = \mathcal{K}P\mathcal{K} + \rho\mathcal{K}$. Let $v_1, \dots, v_N$ be the column vectors of $\mathcal{K}$.

(3) Compute, for $g \in \mathcal{R}_Y$,

$$v_g = \frac{1}{N_g}\sum_{k=1}^{N} v_k \mathbf{1}_{y_k=g}, \quad \bar{v} = \frac{1}{N}\sum_{k=1}^{N} v_k$$

and let $Q = \frac{1}{N}\sum_{g \in \mathcal{R}_Y} N_g(v_g - \bar{v})(v_g - \bar{v})^T$.

(4) Fix $r_0 \le q-1$ and compute the eigenvectors $f_1, \ldots, f_{r_0}$ associated with the $r_0$ largest eigenvalues for the generalized eigenvalue problem $Qf = \lambda Rf$, normalized such that $f_j^T R f_j = 1$.

(5) Compute the scores $\gamma_{jg} = (v_g - \bar{v})^T f_j$.

(6) Given a new observation $x$, let $v(x)$ be the vector with coordinates $K(x, x_k)$, $k = 1, \ldots, N$. Compute the scores $\gamma_j(x) = (v(x) - \bar{v})^T f_j$, $j = 1, \ldots, r_0$. Classify $x$ in the class $g$ maximizing

$$\sum_{i=1}^{r} \gamma_i(x)\gamma_{ig} - \frac{1}{2}\sum_{i=1}^{r} \gamma_{ig}^2 + \log \pi_g. \tag{8.17}$$

## 8.3 Optimal Scoring

It is possible to apply linear regression (chapter 7) to solve a classification problem by mapping the set $\mathcal{R}_Y$ to a collection of $r$-dimensional row vectors, or "scores." These scores (which have a different meaning from the LDA scores) will be represented by a function $\theta : \mathcal{R}_Y \mapsto \mathbb{R}^r$. As an example, one can take $r = q$ and

$$\theta(g_1) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \ \theta(g_2) = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \ldots, \ \theta(g_q) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}.$$

Given a training set $T = (x_1, y_1, \ldots, x_N, y_N)$ and a score function $\theta$, a linear model can then be estimated from data by minimizing

$$\sum_{k=1}^{N} |\theta_{y_k} - a_0 - b^T x_k|^2$$

where $b$ is a $d \times q$ matrix and $a_0 \in \mathbb{R}^q$. Letting as before $\beta$ be the matrix with $a_0^T$ added as first row to $b$ and $\mathcal{X}$ the matrix with first row containing only ones and subsequent rows given by $x_1^T, \ldots, x_N^T$, one gets the least square estimator $\hat{\beta} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathcal{Y}$, where $\mathcal{Y}$ is the $N \times q$ matrix of stacked $\theta_{y_k}^T$ row vectors.

Given an input vector $x$, the row vector $\tilde{x}^T \beta$ will generally not coincide with one of the score vectors. Assignment to a class can then be made by minimizing $|a_0 + b^T x - \theta_g|$ over all $g$ in $\mathcal{R}_Y$.

Since the scores $\theta$ are free to choose, one may also try to optimize them, resulting in the *optimal scoring* algorithm. To describe it, we will need the notation already

introduced for LDA, plus the following.  We will write, for short, $\theta_j = \theta(g_j)$ and introduce the $q \times r$ matrix $\Theta = \begin{pmatrix} \theta_1^T \\ \vdots \\ \theta_q^T \end{pmatrix}$. We also denote by $\rho_1, \ldots, \rho_r$ the column vectors of $\Theta$, so that $\Theta = [\rho_1, \ldots, \rho_r]$. Let $u_{g_i}$, for $i = 1, \ldots, q$, denote the $q$-dimensional vector with $i$th coordinate equal to 1 and all others equal to 0. As before, $N_g$ denote the class sizes, $c_g = N_g/N$, $C$ is the diagonal matrix with coefficients $c_{g_1}, \ldots, c_{g_q}$ and $\zeta = \begin{pmatrix} c_{g_1} \\ \vdots \\ c_{g_q} \end{pmatrix}$.

The goal of optimal scoring is to minimize, now with respect to $\theta$, $a_0$ and $b$, the function

$$F(\theta, a_0, b) = \sum_{k=1}^{N} |\theta(y_k) - a_0 - b^T x_k|^2.$$

Some normalizing conditions are clearly needed, because this problem is under-constrained.  (In the form above, the optimal choice is to take all free parameters equal to 0.)  We now discuss the various indeterminacies and redundancies in the model,

(a) If $R$ is an $r \times r$ orthogonal matrix, then $F(R\theta, Ra_0, bR^T) = F(\theta, a_0, b)$, yielding an infinity of possible equivalent solutions (that all lead to the same classification rule). This implies that there is no loss of generality in assuming that $\Theta^T C\Theta$ is diagonal (introducing $C$ here will turn out to be convenient). Indeed, given any $(\theta, a_0, b)$, one can just take $R$ such that $R\Theta^T C\Theta R^T$ is diagonal and replace $\Theta$ by $R\Theta$, $a_0$ by $Ra_0$ and $b$ by $bR^T$ to get an equivalent solution satisfying the constraint.

(b) Let $D$ be an $r$ by $r$ diagonal matrix with positive entries.  Replace $\theta$, $a_0$ and $b$ respectively by $D\theta$, $Da_0$ and $bD$. The resulting objective function is

$$F(D\theta, Da_0, bD^T) = \sum_{k=1}^{N} |D\theta(y_k) - Da_0 - Db^T x_k|^2$$

$$= \sum_{j=1}^{r} \sum_{k=1}^{N} d_{jj}^2 \left( \theta(y_k, j) - a_0(j) - \sum_{i=1}^{d} b(i,j) x_k(i) \right)^2$$

If the coefficient $d_{jj}$ is free to chose, then the objective function can always be reduced by letting $d_{jj} \to 0$, which removes one of the dimensions in $\theta$. In order to avoid this, one needs to fix the diagonal values of $\Theta^T C\Theta$, and, by symmetry, it is natural to require $\Theta^T C\Theta = \text{Id}_{\mathbb{R}^r}$.

(c) Given any $\delta \in \mathbb{R}^r$, one has $F(\theta, a_0, b) = F(\theta - \delta, a_0 + \delta, b)$, with identical classification rule. One can therefore without loss of generality introduce $r$ linear constraints,

and a convenient choice is

$$\Theta^T \zeta = \sum_{g \in \mathcal{R}_Y} c_g \theta_g = 0.$$

Given this reduction, we can now describe the optimal scoring problem as the minimization of

$$\sum_{k=1}^{N} |\theta_{y_k} - a_0 - b^T x_k|^2$$

subject to $\Theta^T C \Theta = \mathrm{Id}_{\mathbb{R}^r}$ and $\Theta^T \zeta = 0$.

The optimal $a_0$ is given by

$$\hat{a_0} = \frac{1}{N} \sum_{k=1}^{N} \theta_{y_k} - b^T \mu = -b^T \mu,$$

so that the problem is reduced to minimizing

$$\sum_{k=1}^{N} |\theta_{y_k} - b^T (x_k - \mu)|^2$$

subject to the same constraints. Using the facts that $\theta_{y_k} = \Theta^T u_{y_k}$, that

$$\sum_{k=1}^{N} u_{y_k} u_{y_k}^T = \sum_{g \in \mathcal{R}_Y} N_g u_g u_g^T = NC$$

and that

$$\sum_{k=1}^{N} u_{y_k} (x_k - \mu)^T = \sum_{g \in \mathcal{R}_Y} u_g \sum_{k:y_k=g} (x_k - \mu)^T = \sum_{g \in \mathcal{R}_Y} u_g N_g (\mu_g - \mu)^T = NCM,$$

one can write

$$\sum_{k=1}^{N} |\theta_{y_k} - b^T (x_k - \mu)|^2 = \sum_{k=1}^{N} |\Theta^T u_{y_k} - b^T (x_k - \mu)|^2$$

$$= \sum_{k=1}^{N} u_{y_k}^T \Theta \Theta^T u_{y_k} - 2 \sum_{k=1}^{N} (x_k - \mu)^T b \Theta^T u_{y_k} + \sum_{k=1}^{N} (x_k - \mu)^T b b^T (x_k - \mu)$$

$$= \sum_{k=1}^{N} \mathrm{trace}(\Theta^T u_{y_k} u_{y_k}^T \Theta) - 2 \sum_{k=1}^{N} \mathrm{trace}(\Theta^T u_{y_k} (x_k - \mu)^T b)$$

$$+ \sum_{k=1}^{N} \mathrm{trace}(b^T (x_k - \mu)(x_k - \mu)^T b)$$

$$= N \mathrm{trace}(\Theta^T C \Theta) - 2N \mathrm{trace}(\Theta^T CMb) + N \mathrm{trace}(b^T \Sigma_{XX} b).$$

Note that, since $\Theta^T C\Theta = \mathrm{Id}_{\mathbb{R}^r}$, then $\mathrm{trace}(\Theta^T C\Theta) = r$. We therefore obtain a concise form of the optimal scoring problem: minimize

$$-2\mathrm{trace}(\Theta^T CMb) + \mathrm{trace}(b^T \Sigma_{XX} b).$$

subject to $\Theta^T C\Theta = \mathrm{Id}_{\mathbb{R}^r}$ and $\Theta^T \zeta = 0$.

Given $\Theta$, the optimal $b$ is $\Sigma_{XX}^{-1} M^T C\Theta$, and replacing it in the objective function, one finds that $\Theta$ must minimize

$$-2\mathrm{trace}(\Theta^T CM\Sigma_{XX}^{-1} M^T C\Theta) + \mathrm{trace}(\Theta^T CM\Sigma_{XX}^{-1} M^T C\Theta)$$

i.e., maximize

$$\mathrm{trace}(\Theta^T CM\Sigma_{XX}^{-1} M^T C\Theta)$$

subject to $\Theta^T C\Theta = \mathrm{Id}_{\mathbb{R}^r}$ and $\Theta^T \zeta = 0$. We now recall the following linear algebra result (see chapter 2).

**Proposition 8.7** *Let A and B be respectively positive definite and non-negative semidefinite symmetric q by q matrices. Then, the maximum, over all q by r matrices S such that* $\mathrm{trace}(S^T AS) = \mathrm{Id}_{\mathbb{R}^r}$, *of* $\mathrm{trace}(S^T BS)$ *is attained at* $S = [\sigma_1, \ldots, \sigma_r]$, *where the columns vectors* $\sigma_1, \ldots, \sigma_r$ *are the solutions of the generalized eigenvalue problem*

$$B\sigma = \lambda A\sigma$$

*associated with the largest eigenvalues, normalized so that* $\sigma_i^T A\sigma_i = 1$ *for* $i = 1, \ldots, r$..

Given this proposition, let $\rho_1, \ldots, \rho_r$ be the $r$ first eigenvectors for the problem

$$CM\Sigma_{XX}^{-1} M^T C\rho = \lambda C\rho. \tag{8.18}$$

Assume that $r$ is small enough so that the associated eigenvalues are not zero. Let $\Theta = [\rho_1, \ldots, \rho_r]$. We now prove that $\Theta$ is indeed a solution of the optimal scoring problem, and the only point to show to complete the statement is that this $\Theta$ satisfies the constraints $\Theta^T \zeta = 0$. But we have

$$M^T C\mathbb{1}_q = \sum_g c_g (\mu_g - \bar{\mu}) = 0,$$

which implies that $\mathbb{1}_q$ is a solution of the generalized eigenvalue problem associated with $\lambda = 0$. This in turn implies that $\mathbb{1}_q^T C\rho_i = \zeta^T \rho_i = 0$, which is exactly $\Theta^T \zeta = 0$.

To summarize, we have found that the solution $\theta, b$ minimizing

$$-2\mathrm{trace}(\Theta^T CMb) + \mathrm{trace}(b^T \Sigma_{XX} b)$$

subject to $\Theta^T C\Theta = \mathrm{Id}_{\mathbb{R}^r}$ and $\Theta^T \zeta = 0$ is given by

(i) $\Theta = [\rho_1, \ldots, \rho_r]$ where $\rho_1, \ldots, \rho_r$ are the eigenvectors for the problem

$$CM\Sigma_{XX}^{-1}M^T C\rho = \lambda C\rho$$

associated with the $r$ largest eigenvalues, normalized so that $\rho^T C\rho = 1$.

(ii) $b = \Sigma_{XX}^{-1}M^T C\Theta$.

The computation can, however, be further simplified. Let $\lambda_1, \ldots, \lambda_r$ be the eigenvalues associated with $\rho_1, \ldots, \rho_r$. Letting $D$ be the associated diagonal matrix, one can write

$$CM\Sigma_{XX}^{-1}M^T C\Theta = C\Theta D.$$

This yields

$$\Theta = M\Sigma_{XX}^{-1}M^T C\Theta D^{-1} = MbD^{-1},$$

from which we deduce that $\theta_g = \Theta^T u_g = D^{-1}b^T(\mu_g - \bar{\mu})$. So, given a new input vector $x$, the decision rule is to assign it to the class $g$ for which

$$|\theta_g - b^T(x - \bar{\mu})|^2 = |\Theta^T u_g - b^T(x - \bar{\mu})|^2 = |D^{-1}b^T(\mu_g - \bar{\mu}) - b^T(x - \bar{\mu})|^2$$

is minimal. Letting $b_1, \ldots, b_r$ denote the $r$ columns of $b$, this is equivalent to minimizing, in $g$

$$\sum_{j=1}^{r}(b_j^T(\mu_g - \bar{\mu}))^2/\lambda_j^2 - 2\sum_{j=1}^{r}(b_j^T(x - \bar{\mu}))(b_j^T(\mu_g - \bar{\mu}))/\lambda_j. \tag{8.19}$$

From $b = \Sigma_{XX}^{-1}M^T C\Theta$ and $\Theta = MbD^{-1}$ we see that

$$bD = \Sigma_{XX}^{-1}M^T CMb,$$

so that $\Sigma_b b = \Sigma_{XX}bD$. This shows that the columns of $b$ are solution of the eigenvalue problem $\Sigma_b u = \lambda \Sigma_{XX}u$. Moreover, from $\Theta^T C\Theta = \mathrm{Id}_{\mathbb{R}^r}$, we get $b^T \Sigma_b b = D^2$. Since $b^T \Sigma_b b = b^T \Sigma_{XX}bD$, we get that $b$ must be normalized to that $b^T \Sigma_{XX}b = D$.

This shows that the solution of the optimal scoring problem can be reformulated uniquely in terms of $b$: if $b_1, \ldots, b_r$ are the $r$ principal solutions of the eigenvalue problem $\Sigma_b u = \lambda \Sigma_{XX}u$, normalized so that $u^T \Sigma_{XX}u = \lambda$, a new input is classified into the class $g$ minimizing

$$\sum_{j=1}^{r}\gamma_j(\mu_g)^2/\lambda_j^2 - 2\sum_{j=1}^{r}\gamma_j(x)\gamma_j(\mu_g)/\lambda_j.$$

with $\gamma_j(x) = b_j^T(x - \bar{\mu})$.

**Remark 8.8** The following computation shows that optimal scoring is closely related to LDA. Recall the identity $\Sigma_{XX} = \Sigma_w + \Sigma_b$. It implies that a solution of $\Sigma_b u = \lambda \Sigma_{XX} u$ is also a solution of $\Sigma_b u = \tilde{\lambda} \Sigma_w u$ with $\tilde{\lambda} = \lambda/(1 - \lambda)$. If $u^T \Sigma_{XX} u = \lambda$, then

$$u^T \Sigma_w u = \lambda - u^T \Sigma_b u = \lambda - \lambda^2 = \frac{\tilde{\lambda}}{(1 + \tilde{\lambda})^2},$$

which shows that

$$\tilde{u} = \frac{1 + \tilde{\lambda}}{\sqrt{\tilde{\lambda}}} u$$

satisfies $\tilde{u}^T \Sigma_w \tilde{u} = 1$. So,

$$e_j = \frac{1 + \tilde{\lambda}_j}{\sqrt{\tilde{\lambda}_j}} b_j$$

coincide with the LDA directions. We have, letting $\tilde{\gamma}_j(x) = e_j^T(x - \bar{\mu}) = \sqrt{\tilde{\lambda}_j}\gamma_j(x)/(1 + \tilde{\lambda}_j)$:

$$\sum_{j=1}^{r} \gamma_j(\mu_g)^2/\lambda_j^2 - 2\sum_{j=1}^{r} \gamma_j(x)\gamma_j(\mu_g)/\lambda_j = \sum_{j=1}^{r} \tilde{\gamma}_j(\mu_g)^2/\tilde{\lambda}_j - 2\sum_{j=1}^{r} \gamma_j(x)\gamma_j(\mu_g)/(1 + \tilde{\lambda}_j)$$

which relates the classification rules for the two methods.                    ◆

**Remark 8.9** Optimal scoring can be modified by adding a penalty in the form

$$\gamma \sum_{i=1}^{r} b_i^T \Omega b_i = \gamma \text{trace}(b^T \Omega b) \tag{8.20}$$

where $\Omega$ is a weight matrix. This only modifies the previous discussion by adding $\gamma\Omega/N$ to both $\Sigma_{XX}$ and $\Sigma_w$.                    ◆

### 8.3.1   Kernel optimal scoring

Let $h : \mathcal{R}_X \to H$ be the feature function and $K$ the associated kernel, as usual. Optimal scoring in feature space requires to minimize

$$\sum_{k=1}^{N} |\theta_{y_k} - a_0 - \mathbb{b}(h(x_k))|^2 + \gamma \|\mathbb{b}\|_H^2,$$

where we have introduced a penalty on $\mathbb{b}$. Here, $\mathbb{b}$ is a linear operator from $H$ to $\mathbb{R}^r$, therefore taking the form

$$\mathbb{b}(h) = \begin{pmatrix} \langle b_1, h \rangle_H \\ \vdots \\ \langle b_r, h \rangle_H \end{pmatrix}$$

with $b_1, \ldots, b_r \in H$, and we take

$$\|\mathbb{b}\|_H^2 = \sum_{i=1}^r \|b_i\|_H^2.$$

It is once again clear (and the argument is left to the reader) that the problem can be reduced to the finite dimensional space $V = \mathrm{span}(h(x_1), \ldots, h(x_N))$, and that the optimal $b_1, \ldots, b_r$ must take the form

$$b_j = \sum_{l=1}^N \alpha_{li} h(x_l).$$

Introduce the kernel matrix $\mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$ with $k$th column denoted $\mathcal{K}^{(k)}$. Let $\alpha$ be the $N$ by $r$ matrix with entries $\alpha_{kj}$, $k = 1, \ldots, N$, $j = 1, \ldots, r$. Then $\mathbb{b}(h(x_k))$, which is the vector with coordinates

$$\langle b_j, h(x_k) \rangle = \sum_{l=1}^N \alpha_{li} K(x_k, x_l),$$

is equal to $\alpha^T \mathcal{K}^{(k)}$. Moreover

$$\|\mathbb{b}\|_H^2 = \sum_{j=1}^r \sum_{k,l=1}^N \alpha_{kj} K(x_k, x_l) \alpha_{lj} = \mathrm{trace}(\alpha^T \mathcal{K} \alpha).$$

We therefore need to minimize

$$\sum_{k=1}^N |\theta_{y_k} - a_0 - \alpha^T \mathcal{K}^{(k)}|^2 + \gamma \, \mathrm{trace}(\alpha^T \mathcal{K} \alpha),$$

so that the problem is reduced to penalized optimal scoring, with $x_k$ replaced by $\mathcal{K}^{(k)}$, $b$ replaced by $\alpha$ and the matrix $\Omega$ in (8.20) replaced by $\mathcal{K}$. Introducing the matrix $P = \mathrm{Id}_{\mathbb{R}^N} - \mathbb{1}\mathbb{1}^T/N$ and $\mathcal{K}_c = P\mathcal{K}$, the covariance matrix $\Sigma_{XX}$ becomes $\mathcal{K}_c^T \mathcal{K}_c/N = \mathcal{K} P \mathcal{K}/N$.

The class averages $\mu_g$ are equal to $\mathcal{K}\mathbb{1}(g)/N_g$ while $\mu = \mathcal{K}\mathbb{1}/N$, so that the matrix $M$ is equal to

$$\begin{pmatrix} \mathbb{1}(g_1)^T/N_{g_1} - \mathbb{1}^T/N \\ \vdots \\ \mathbb{1}(g_q)^T/N_{g_q} - \mathbb{1}^T/N \end{pmatrix} \mathcal{K}$$

which gives $\Sigma_b = M^T C M = \mathcal{K} Q \mathcal{K}$, where $Q$ is

$$Q = PCP = \sum_{g \in \mathcal{R}_Y} \frac{N_g}{N} \left( \frac{\mathbb{1}(g)}{N_g} - \frac{\mathbb{1}}{N} \right) \left( \frac{\mathbb{1}(g)}{N_g} - \frac{\mathbb{1}}{N} \right)^T$$

So, the columns of $\alpha$ are the $r$ principal eigenvectors $\rho_1, \ldots, \rho_r$ of the problem

$$\mathcal{K} Q \mathcal{K} \rho = \frac{1}{N} (\mathcal{K} P \mathcal{K} + \gamma \mathcal{K}) \rho.$$

Given $\alpha$, one then has, for any $x \in \mathbb{R}^d$,

$$\langle b_i, h(x) \rangle_H = \sum_{k=1}^N \alpha_{ki} K(x, x_k)$$

and

$$a_0^{(i)} = \frac{1}{N} \sum_{k,l=1}^N \alpha_{ki} K(x_l, x_k).$$

## 8.4   Separating hyperplanes and SVMs

### 8.4.1   One-layer perceptron and margin

In this whole section, we restrict to two-class problems, and let $\mathcal{R}_Y = \{-1, 1\}$. Given $a_0 \in \mathbb{R}$ and $b \neq 0 \in \mathbb{R}^d$, the equation $a_0 + b^T x = 0$ defines a hyperplane in $\mathbb{R}^d$. The function $f(x) = \text{sign}(a_0 + x^T b)$ defines a classifier that attributes a class $\pm 1$ to $x$ according to which side of the hyperplane it belongs (we ignore the ambiguity when $x$ is on the hyperplane). With this notation, a pair $(x, y)$, where $y$ is the true class, is correctly classified if and only if $y(a_0 + x^T b) > 0$.

Let $T = (x_1, y_1), \ldots, (x_N, y_N)$ denote, as usual, the training data. A hyperplane, represented by the parameters $(a_0, b)$ is *separating* for $T$ if it correctly classifies all its samples, i.e., if $y_k(a_0 + x_k^T b) > 0$ for $k = 1, \ldots, N$. If such a hyperplane exists, one says that $T$ is linearly separable.

The perceptron algorithm computes $a_0$ and $b$ by minimizing

$$L(\beta) = \sum_{k=1}^N [y_k(a_0 + x_k^T b)]^-$$

with $u^- = \max(0, -u)$, , or more precisely, fixing a small positive number $\delta$:

$$L(\beta) = \sum_{k=1}^{N} [\delta - y_k(a_0 + x_k^T b)]^+ .$$

The problem can be recast as a linear program, i.e., minimize

$$\sum_{k=1}^{N} \xi_k$$

subject to $\xi_k \geq 0, \xi_k + y_k(a_0 + x_k^T b) - \delta \geq 0$ for $k = 1, \ldots, N$.

However, when $T$ is linearly separable, separating hyperplanes are not uniquely defined, and there is in general (depending on the choice made for $\delta$) an infinity of solutions to the perceptron problem. Intuitively, one should prefers a solution that classifies the training data with some large margin, rather than one for which training points may be very close to the separating boundary (see fig. 8.2).



Figure 8.2: The green line is preferable to the purple one in order to separate the data.

This leads to the maximum margin separating hyperplane classifier, also called linear SVM, introduced by Vapnik and Chervonenkis [196, 197].

### 8.4.2 Maximizing the margin

We will use the following result.

**Proposition 8.10** *The distance of a point $x \in \mathbb{R}^d$ to the hyperplane $M : a_0 + b^T x = 0$ is given by $|a_0 + x^T b|/|b|$.*

PROOF By definition, $\text{distance}(x, M) = |x - \pi_M(x)|$ where $\pi_M$ is the orthogonal projection on $M$. Since $b$ is normal to $M$ and letting $h = \pi_M(x)$, we have $x = \lambda b + h$ so that $|\lambda b| = \text{distance}(x, M)$. Writing $a_0 + b^T h = 0$ in this equation implies $a_0 + b^T x = \lambda |b|^2$ so that $|\lambda||b| = |a_0 + x^T b|/|b|$. ∎

Assume that $T$ is linearly separable and let $M : a_0 + b^T x = 0$ be a separating hyperplane. The classification margin is defined as the minimal distance of the input vectors $x_1, \ldots, x_N$ to this hyperplane, i.e.,

$$m(a_0, b) = \min\{|a_0 + x_k^T b|/|b| : k = 1, \ldots, N\}.$$

Because the hyperplane is separating, we have $y_k(a_0 + x_k^T b) = |a_0 + x_k^T b|$ for all $k$, so that we also have

$$m(a_0, b) = \min\{y_k(a_0 + x_k^T b)/|b| : k = 1, \ldots, N\}.$$

We want to maximize this margin among all separating hyperplanes. This can be expressed as maximizing, with respect to $a_0, b$, the quantity

$$\min\{y_k(a_0 + x_k^T b)/|b| : k = 1, \ldots, N\}$$

subject to the constraint that the hyperplane is separating, namely

$$y_k(a_0 + x_k^T b) \geq 0, \ k = 1, \ldots, N.$$

Introducing a new variable $C$ representing the margin, the previous problem is equivalent to maximizing $C$ subject to

$$y_k(a_0 + x_k^T b) \geq C|b|, \ k = 1, \ldots, N.$$

The problem is now overparametrized, and there is no loss of generality in enforcing the additional constraint $C|b| = 1$. Noting that maximizing $C$ is the same as minimizing $|b|^2$, we can now reformulate the maximum margin hyperplane problem as minimizing $|b|^2/2$ subject to

$$y_k(a_0 + x_k^T b) \geq 1, \ k = 1, \ldots, N,$$

with $C$ (the margin) given by $C = 1/|b|$. This results in a quadratic programming problem.

If the data is not separable, there is no feasible point for this problem. To also account for this situation (which is common), we can replace the constraint by a penalty and minimize, with respect to $a_0$ and $b$:

$$\frac{|b|^2}{2} + \gamma \sum_{k=1}^{N} (1 - y_k(a_0 + x_k^T b))^+$$

for some $\gamma > 0$. (Recall that $x^+ = \max(x, 0)$.) This is equivalent to minimizing the perceptron objective function, with $\delta = 1$, and with an additional penalty term equal to $|b|^2/(2\gamma)$. This minimization problem is equivalent to a quadratic programming problem obtained by introducing slack variables $\xi_k$, $k = 1, \ldots, N$ and minimizing

$$\frac{1}{2}|b|^2 + \gamma \sum_{k=1}^{N} \xi_k,$$

subject to the constraints $\xi_k \geq 0$, $y_k(a_0 + x_k^T b) + \xi_k \geq 1$, for $k = 1, \ldots, N$.

### 8.4.3 KKT conditions and dual problem

Introduce Lagrange multipliers $\eta_k \geq 0$ for $\xi_k \geq 0$ and $\alpha_k \geq 0$ for $y_k(a_0 + x_k^T b) + \xi_k \geq 1$. The Lagrangian is then given by

$$\mathcal{L} = \frac{1}{2}|b|^2 + \gamma \sum_{k=1}^{N} \xi_k - \sum_{k=1}^{N} \eta_k \xi_k - \sum_{k=1}^{N} \alpha_k \left( y_k(a_0 + x_k^T b) + \xi_k - 1 \right).$$

The KKT conditions are

$$\begin{cases} b - \sum_{k=1}^{N} \alpha_k y_k x_k = 0 \\[2mm] \sum_{k=1}^{N} \alpha_k y_k = 0 \\[2mm] \gamma - \eta_k - \alpha_k = 0, \quad k = 1, \ldots, N \\[1mm] \xi_k \eta_k = 0, \quad k = 1, \ldots, N \\[1mm] \alpha_k \left( y_k(a_0 + x_k^T b) + \xi_k - 1 \right) = 0, \quad k = 1, \ldots, N \end{cases} \tag{8.21}$$

Minimizing $\mathcal{L}$ with respect to $a_0$, $b$ and $\xi_1, \ldots, \xi_N$ and ensuring that the minimum is finite provides the first three KKT conditions. The resulting dual formulation therefore requires to maximize

$$\sum_{k=1}^{N} \alpha_k - \frac{1}{2} \sum_{k,l=1}^{N} \alpha_k \alpha_l y_k y_l x_k^T x_l$$

subject to the constraints $0 \leq \alpha_k \leq \gamma$, $\sum_{k=1}^{N} \alpha_k y_k = 0$.

We now discuss the consequences of the complementary slackness conditions based on the position of training sample relative to the separating hyperplane.

(i) First consider indices $k$ such that $(x_k, y_k)$ is correctly classified *beyond the margin*, i.e., $y_k(a_0 + x_k^T b) > 1$. The last KKT condition and the constraint $\xi_k \geq 0$ require $\alpha_k = 0$, and the third one then gives $\xi_k = 0$.

(ii) For samples that are misclassified or correctly classified below the margin [5], i,e., $y_k(a_0 + x_k^T b) < 1$, the constraint $y_k(a_0 + x_k^T b) + \xi_k \geq 1$ implies $\xi_k > 0$, so that $\alpha_k = \gamma$ and $y_k(a_0 + x_k^T b) + \xi_k = 1$.

(iii) If $(x_k, y_k)$ is correctly classified exactly at the margin, then $\xi_k = 0$ and there is no constrain on $\alpha_k$ beside belonging to $[0, \gamma]$. Training samples that lie exactly at the margin are called *support vectors*.

Given a solution $\alpha_1, \ldots, \alpha_N$ of the dual problem, one immediately recovers $b$ via the first equation in (8.21). For $a_0$, one must, similarly to the regression case, rely on support vectors, which can be identified when $0 < \alpha_k < \gamma$. In this case, one can take $a_0 = y_k - x_k^T b$.

If no support vector is found, then $a_0$ is not uniquely determined, and can be any value such that $y_k(a_0 + b^T x_k) \geq 1$ if $\alpha_k = 0$ and $y_k(a_0 + b^T x_k) \leq 1$ if $\alpha_k = \gamma$. This shows that $a_0$ can be any point in the interval $[\beta_0^-, \beta_0^+]$ with

$$a_0^- = \max\{y_k - x_k^T b : (y_k = 1 \text{ and } \alpha_k = 0) \text{ or } (y_k = -1 \text{ and } \alpha_k = \gamma)\}$$
$$a_0^+ = \min\{y_k - x_k^T b : (y_k = -1 \text{ and } \alpha_k = 0) \text{ or } (y_k = 1 \text{ and } \alpha_k = \gamma)\}.$$

### 8.4.4  Kernel version

We make the usual assumptions: $h : \mathcal{R}_X \to H$ is a feature map with values in an inner-product space with $K(x, y) = \langle h(x), h(y) \rangle_H$. The predictors take the form $f(x) = \text{sign}(a_0 + \langle b, h(x) \rangle_H)$, $a_0 \in \mathbb{R}$ and $b \in H$, and the goal is to minimize

$$\frac{1}{2}\|b\|_H^2 + \gamma \sum_{k=1}^N \xi_k,$$

subject to $\xi_k \geq 0$, $y_k(a_0 + \langle h(x_k), b \rangle_H) + \xi_k \geq 1$, $k = 1, \ldots, N$.

Let $V = \text{span}(h(x_1), \ldots, h(x_N))$. The usual projection argument implies that the optimal $b$ must belong to $V$ and therefore take the form

$$b = \sum_{k=1}^N u_k h(x_k).$$

---

[5]Note that, even if the training data is linearly separable, there are generally samples that are on the right side of the hyperplane, but at a distance to the hyperplane strictly lower that the "nominal margin" $C = 1/|b|$. This is due to our relaxation of the original problem of finding a separating hyperplane with maximal margin.

We therefore need to minimize

$$\frac{1}{2} \sum_{k,l=1}^{N} u_k u_l K(x_k, x_l) + \gamma \sum_{k=1}^{N} \xi_k,$$

subject to

$$y_k \left( a_0 + \sum_{l=1}^{N} K(x_k, x_l) a_l \right) + \xi_k \geq 1$$

for $k = 1, \ldots, N$. Introducing the same Lagrange multipliers as before, the Lagrangian is

$$\mathcal{L} = \frac{1}{2} \sum_{k,l=1}^{N} u_k u_l K(x_k, x_l) + \gamma \sum_{k=1}^{N} \xi_k$$

$$- \sum_{k=1}^{N} \eta_k \xi_k - \sum_{k=1}^{N} \alpha_k \left( y_k \left( a_0 + \sum_{l=1}^{N} K(x_k, x_l) u_l \right) + \xi_k - 1 \right).$$

Using vector notation, we have

$$\mathcal{L} = \frac{1}{2} u^T \mathcal{K} u + \xi^T (\gamma \mathbb{1} - \eta - \alpha) - a_0 \alpha^T y - (\alpha \odot y)^T \mathcal{K} u + \alpha^T \mathbb{1}$$

where $y \odot \alpha$ is the vector with coordinates $y_k \alpha_k$. The infimum of $\mathcal{L}$ is $-\infty$ unless $\gamma \mathbb{1} - \eta - \alpha = 0$ and $\alpha^T y = 0$. If these identities are true, then the optimal $u$ is $u = \alpha \odot y$ and the minimum of $\mathcal{L}$ is

$$-\frac{1}{2} (\alpha \odot y)^T \mathcal{K} (\alpha \odot y) + \alpha^T \mathbb{1}$$

The dual problem therefore requires to minimize

$$\frac{1}{2} (\alpha \odot y)^T \mathcal{K} (\alpha \odot y) - \alpha^T \mathbb{1} = \alpha^T (\mathcal{K} \odot yy^T) \alpha - \alpha^T \mathbb{1}$$

subject to $\gamma \mathbb{1} - \eta - \alpha = 0$ and $\alpha^T y = 0$.

This is exactly the same problem as the one we obtained in the linear case, up to the replacement of the Euclidean inner products $x_k^T x_l$ by the kernel evaluations $K(x_k, x_l)$. Given the solution of the dual problem, the optimal $b$ is

$$b = \sum_{k} u_k h(x_k) = \sum_{k=1}^{N} \alpha_k y_k h(x_k).$$

It is no computable, but the classification rule is explicit and given by

$$f(x) = \text{sign}\left(a_0 + \sum_{k=1}^{N} \alpha_k y_k K(x_k, x)\right).$$

Similarly to the linear case, the coefficient $a_0$ can be identified using a support vector, or is otherwise not uniquely determined. More precisely, if one of the $\alpha_k$'s is strictly between 0 and $\gamma$, then $a_0$ is given by $a_0 = y_k - \sum_l \alpha_l y_l K(x_k, x_l)$. Otherwise, $a_0$ is any number between

$$\beta_0^- = \max\left\{y_k - \sum_l \alpha_l y_l K(x_k, x_l) : (y_k = 1 \text{ and } \alpha_k = 0) \text{ or } (y_k = -1 \text{ and } \alpha_k = \gamma)\right\}$$

and

$$\beta_0^+ = \min\left\{y_k - \sum_l \alpha_l y_l K(x_k, x_l) : (y_k = -1 \text{ and } \alpha_k = 0) \text{ or } (y_k = 1 \text{ and } \alpha_k = \gamma)\right\}.$$

**Chapter 9**

# Nearest-Neighbor Methods

Unlike linear models, nearest-neighbor methods are completely non-parametric and assume no regularity on the decision rule or the regression function. In their simplest version, they require no training and rely on the proximity of a new observation to those that belong to the training set. We will discuss in this chapter how these methods are used for regression and classification, and study some of their theoretical properties.

## 9.1 Nearest neighbors for regression

### 9.1.1 Consistency

We let $\mathcal{R}_X$ denote the input space, and $\mathcal{R}_Y = \mathbb{R}^q$ be the output space. We assume that a distance, denoted dist is defined on $\mathcal{R}_X$. This means that dist $: \mathcal{R}_X \times \mathcal{R}_X \to [0, +\infty]$ (we allow for infinite values) is a symmetric function such that $\text{dist}(x, x') = 0$ if and only if $x = x'$ and, for all $x, x', x'' \in \mathcal{R}_X$

$$\text{dist}(x, x') \leq \text{dist}(x, x'') + \text{dist}(x'', x'),$$

which is the triangle inequality.

Let $T = (x_1, y_1, \ldots, x_N, y_N)$ be the training set. For $x \in \mathcal{R}_X$, let

$$D_T(x) = (\text{dist}(x, x_k), k = 1, \ldots, N)$$

be the collection of all distances between $x$ and the inputs in the training set. We consider regression estimators taking the form

$$\hat{f}(x) = \sum_{k=1}^{N} W_k(x) y_k \tag{9.1}$$

where $W_1(x), \ldots, W_N(x)$ is a family of coefficients, or weights, that only depends on $D_T(x)$.

We will, more precisely, use the following construction [183]. Assume that a family of numbers $w_1 \geq w_2 \geq \cdots \geq w_N \geq 0$ is chosen, with $\sum_{j=1}^N w_j = 1$. Given $x \in \mathbb{R}^d$ and $k \in \{1, \ldots, N\}$, we let $r_k^+(x)$ denote the number of indexes $k'$ such that $\text{dist}(x, x_{k'}) \leq \text{dist}(x, x_k)$ and $r_k^-(x)$ the number of such indexes such that $d(x, x_{k'}) < d(x, x_k)$. The coefficients defining $\hat{f}$ in (9.1) are then chosen as:

$$W_k(x) = \frac{\sum_{k'=r_k^-(x)+1}^{r_k^+(x)} w_{k'}}{r_k^+(x) - r_k^-(x)}. \tag{9.2}$$

To emphasize the role of $(w_1, \ldots, w_N)$ is this definition, we will denote the resulting estimation as $\hat{f}_w$. If there is no tie in the sequence of distances between $x$ and elements of the training set, then $r_k^+(x) = r_k^-(x) + 1$ is the rank of $x_k$ when training data is ordered according to their proximity to $x$, and $W_k(x) = w_{r_k^+(x)}$. In this case, defining $l_1, \ldots, l_N$ such that $d(x, x_{l_1}) < \cdots < d(x, x_{l_N})$, we have

$$\hat{f}_w(x) = \sum_{j=1}^N w_j y_{l_j}.$$

In the general case, the weights $w_j$ associated with tied observations are averaged.

If $p$ is an integer, the $p$-nearest neighbor ($p$-NN) estimator (that we will denote $\hat{f}_p$) is associated to the weights $w_j = 1/p$ for $j = 1, \ldots, p$ and 0 otherwise. If there is no tie for the definition of the $p$th nearest neighbor of $x$, $W_k(x) = 1/p$ if $k$ is among the $p$ nearest-neighbors and $W_k(x) = 0$ otherwise, so that $\hat{f}_p$ is the average of the output values over these $p$ nearest neighbors. If the $p$th nearest neighbors are tied, their output value is averaged before being used in the sum. For example, assume that $N = 5$ and $p = 2$ and let the distances between $x$ and $x_k$ for $k = 1, \ldots, 5$ be respectively $9, 3, 2, 4, 6$. Then $\hat{f}_2(x) = (y_2 + y_3)/2$. If the distances were $9, 3, 2, 3, 6$, then we would have $\hat{f}_2(x) = (y_2 + y_4)/4 + y_3/2$.

When $\mathcal{R}_X = \mathbb{R}^d$ and $d(x, x') = |x - x'|$, the following result is true.

**Theorem 9.1 ([183])** *Assume that $\mathbb{E}(Y^2) < \infty$. Assume that, for each $N$, a sequence $w^{(N)} = w_1^{(N)} \geq \cdots \geq w_N^{(N)} \geq 0$ is chosen with $\sum_{j=1}^N w_j^{(N)} = 1$. Assume, in addition, that*

*(i)* $\lim_{N \to \infty} w_1^{(N)} = 0$

*(ii)* $\lim_{N \to \infty} \sum_{j \geq \alpha N} w_j^{(N)} \to 0$, *for some $\alpha \in (0, 1)$.*

*Then the corresponding classifier $\hat{f}_{w^{(N)}}$ converges in the $L^2$ norm to $\mathbb{E}(Y \mid X)$:*

$$\mathbb{E}\left(|\hat{f}_{w^{(N)}}(X) - \mathbb{E}(Y \mid X)|^2\right) \to 0.$$

*For nearest-neighbor regression, (i) and (ii) mean that the number of nearest neighbors $p_N$ must be chosen such that $p_N \to \infty$ and $p_N/N \to 0$.*

PROOF We give a proof under the assumption that $f : x \mapsto \mathbb{E}(Y \mid X = x)$ is uniformly continuous and bounded (one can, in fact, prove that it is always possible to reduce to this case).

To lighten the notation, we will not make explicit the dependency on $N$ in of quantities such as $W$ or $w$. One has

$$\hat{f}_w(X) - \mathbb{E}(Y \mid X) = \sum_{k=1}^{N} W_k(X)(f(X_k) - f(X)) + \sum_{k=1}^{N} W_k(X)(Y_k - f(X_k)) \qquad (9.3)$$

and the two sums can be addressed separately.

We start with the first sum and write, by Schwartz's inequality:

$$\left(\sum_k W_k(X)(f(X_k) - f(X))\right)^2 \leq \sum_k W_k(X)(f(X_k) - f(X))^2.$$

It therefore suffices to study the limit of $E(\sum_k W_k(X)(f(X_k) - f(X))^2$. Fix $\epsilon > 0$. By assumption, there exists $M, a > 0$ such that $|f(x)| \leq M$ for all $x$ and $|x - y| \leq a \Rightarrow |f(x) - f(y)|^2 \leq \epsilon$. Then

$$\mathbb{E}\left(\sum_k W_k(X)(f(X_k) - f(X))^2\right) = \mathbb{E}\left(\sum_k W_k(X)(f(X_k) - f(X))^2 \, \mathbf{1}_{|X_k - X| \leq a}\right)$$

$$+ \mathbb{E}\left(\sum_k W_k(X)(f(X_k) - f(X))^2 \, \mathbf{1}_{|X_k - X| > a}\right)$$

$$\leq \epsilon^2 + 4M^2 \mathbb{E}\left(\sum_k W_k(X)\mathbf{1}_{|X_k - X| > a}\right).$$

Since $\epsilon$ can be made arbitrarily small, we need to show that, for any positive $a$, the second term in the upper-bound tends to 0 when $N \to \infty$. We will use the following fact, which requires some minor measure theory argument to prove rigorously. Define

$$S = \{x : \forall \delta > 0, \mathbb{P}(|X - x| < \delta) > 0\}.$$

This set is called the support of $X$. Then, one can show that $\mathbb{P}(X \in S) = 1$. This means that, if $\tilde{X}$ is independent from $X$ with the same distribution, then, for any $\delta > 0$, $\mathbb{P}(|X - \tilde{X}| < \delta | X) > 0$ with probability one. [1]

Let $N_a(x) = |\{k : |X_k - x| \leq a\}|$. We have, for all $x \in S$ and $a > 0$, and using the law of large numbers,

$$\frac{N_a(x)}{N} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{1}_{|X_k - x| \leq a} \to P(|X - x| \leq a) > 0.$$

If $|X - X_k| > a$, then $r_k^-(X) > N_a(x)$ so that

$$\sum_k W_k(X) \mathbf{1}_{|X_k - X| > a} \leq \sum_{j \geq N_a(X)} w_j,$$

and we have, taking $0 < \alpha < P(|X - x| \leq a)$,

$$\mathbb{E}\left( \sum_{j \geq N_a(X)} w_j \right) \leq \sum_{j \geq \alpha N} w_j + \mathbb{P}(N_a(X) < \alpha N)$$

and both terms in the upper bound converge to 0. This shows that the first sum in (9.3) tends to 0.

We now consider the second sum in (9.3). Let $Z_k = Y_k - \mathbb{E}(Y \mid X_k)$. We have $\mathbb{E}(Z_k \mid X_k) = 0$ and $\mathbb{E}(Z_k^2) < \infty$. We can write

$$\mathbb{E}\left( \left| \sum_{k=1}^{N} W_k(X) Z_k \right|^2 \right) = \mathbb{E}\left( \mathbb{E}\left( \left| \sum_{k=1}^{N} W_k(X) Z_k \right|^2 \Bigg| X, X_1, \ldots, X_N \right) \right)$$

$$= \mathbb{E}\left( \sum_{k=1}^{N} W_k(X)^2 \mathbb{E}(Z_k^2 \mid X_k) \right)$$

$$+ \sum_{k \neq l = 1}^{N} \mathbb{E}(W_k(X) W_l(X) \mathbb{E}(Z_k Z_l \mid X_i, X_j))$$

---

[1] This statement is proved as follows (with the assumption that $X$ is Borel measurable). Let $S^c$ denote the complement of $S$. Then $S^c$ is open. Indeed if $x \notin S$, there exists $\delta_x > 0$ such that, letting $B(x, \delta_x)$ denote the open ball with radius $\delta_x$, $\mathbb{P}(X \in B(x, \delta_x)) = 0$. Then $\mathbb{P}(X \in B(x', \delta_x/3)) = 0$ as soon as $|x - x'| < \delta_x/3$, so that $B(x, \delta_x/3) \subset S^c$.

If $K \subset S^c$ is compact, then $K \subset \bigcup_{x \in K} B(x, \delta_x)$ and one can find a finite subset $M \subset K$ such that $K \subset \bigcup_{x \in M} B(x, \delta_x)$, which proves that $\mathbb{P}(X \in K) = 0$. Since $\mathbb{P}(X \in S^c) = \max_K \mathbb{P}(X \in K)$ where the maximum is over all compact subsets of $S^c$, we find $\mathbb{P}(X \in S^c) = 0$ as required.

The cross products in the last term vanish because $\mathbb{E}(Z_k \mid X_k) = 0$ and the samples are independent. So it only remains to consider

$$\mathbb{E}\left(\sum_{k=1}^{N} W_k(X)^2 \mathbb{E}(Z_k^2 \mid X_k)\right)$$

The random variable $\mathbb{E}(Z_k \mid X_k) = \mathbb{E}(Y_k^2 \mid X_k) - \mathbb{E}(Y_k \mid X_k)^2$ is a fixed non-negative function of $X_k$, that we will denote $h(X_k)$. We have

$$\mathbb{E}\left(\sum_{k=1}^{N} W_k(X)^2 h(X_k)\right) \leq w_1 \mathbb{E}\left(\sum_{k=1}^{N} W_k(X) h(X_i)\right)$$

with $w_1 \to 0$ and the proof is concluded by showing that $\mathbb{E}\left(\sum_{k=1}^{N} W_k(X) h(X_k)\right)$ is bounded.

Recall that the weights $W_k$ are functions of $X$ and of the whole training set, and we will need to make this dependency explicit and write $W_i(X, \mathbb{T}_X)$ where $\mathbb{T}_X = (X_1, \ldots, X_N)$. Similarly, the ranks in (9.2) will be written $r_j^+(X, \mathbb{T}_X)$ and $r_j^-(X, \mathbb{T}_X)$.

Because $X, X_1, \ldots, X_N$ are i.i.d., we can switch the role of $X$ and $X_k$ in the $k$th term of the sum, yielding

$$\mathbb{E}\left(\sum_{k=1}^{N} W_k(X, \mathbb{T}_X) h(X_k)\right) = \mathbb{E}\left(\left(\sum_{i=1}^{N} W_k(X_k, \mathbb{T}_X^{(k)})\right) h(X)\right)$$

with $\mathbb{T}_X^{(k)} = (X_1, \ldots, X_{k-1}, X, X_{k+1}, \ldots, X_N)$. We now show that $\sum_{k=1}^{N} W_k(X_k, \mathbb{T}_X^{(k)})$ is bounded independently of $X, X_1, \ldots, X_N$.

For this purpose, we group $X_1, \ldots, X_N$ according to approximate alignment with $X$. For $u \in \mathbb{R}^d$ with $|u| = 1$ and for $\delta \in (0, \pi/4)$, denote by $\Gamma(u, \delta)$ the cone formed by all vectors $v$ in $\mathbb{R}^d$ such that $\langle v, u \rangle > |v| \cos \delta$ (i.e., the angle between $v$ and $u$ is less than $\delta$). Notice that if $v, v' \in \Gamma(u, \delta)$, then $\langle v, v' \rangle \geq \cos(2\delta)|v||v'|$ and if $|v'| \leq |v|$, then

$$|v|^2 - |v - v'|^2 = |v'|(2|v|\cos(2\delta) - |v'|) > 0 \tag{9.4}$$

because $\cos(2\delta) > 1/2$.

Fixing $\delta$, let $C_d(\delta)$ be the minimal number of such cones needed to cover $\mathbb{R}^d$. Choosing such a covering $\Gamma(u_1, \delta), \ldots, \Gamma(u_M, \delta)$ where $M = C_d(\delta)$, we define the following subsets of $\{1, \ldots, M\}$:

$$I_0 = \{k : X_k = X\}$$
$$I_q = \left\{k \notin I_0 : X_k - X \in \Gamma(u_q, \delta)\right\}, \quad q = 1, \ldots, M$$

(these sets may overlap). We have

$$\sum_{k=1}^{N} W_k(X_k, \mathbb{T}_X^{(k)}) \le \sum_{q=0}^{M} \sum_{k \in I_q} W_k(X_k, \mathbb{T}_X^{(k)})$$

If $k \in I_0$, then $r_k^-(X_k, \mathbb{T}_X^{(k)}) = 0$ and $r_k^+(X_k, \mathbb{T}_X^{(k)}) = c$ with $c = |I_0|$. This implies that, for $k \in I_0$, we have $W_k(X_k, \mathbb{T}_X^{(k)}) = \sum_{j=1}^{c} w_j/c$ and

$$\sum_{k \in I_0} W_k(X_k, \mathbb{T}_X^{(k)}) = \sum_{j=1}^{c} w_j.$$

We now consider $I_q$ with $q \ge 1$. Write $I_q = \{i_1, \dots, i_r\}$ ordered so that $|X_{i_j} - X|$ is non-decreasing. If $j' < j$, we have (using (9.4)) $|X_{i_j} - X_{i_{j'}}| < |X - X_{i_j}|$. This implies that $r_{i_j}^-(X_{i_j}, \mathbb{T}_X^{(i_j)}) \ge j - 1$ and $r_{i_j}^+(X_{i_j}, \mathbb{T}_X^{(i_j)}) - r^- i_j(X_{i_j}, \mathbb{T}_X^{(i_j)}) \ge c + 1$. Therefore,

$$W_{i_j}(X_{i_j}, \mathbb{T}_X^{(i_j)}) \le \frac{1}{c+1} \sum_{j'=j}^{c+j} w_{j'}$$

and

$$\sum_{k \in I_q} W_k(X_k, \mathbb{T}_X^{(k)}) \le \frac{1}{c+1} \sum_{j=1}^{N} \sum_{j'=j}^{c+j} w_{j'} = \frac{1}{c+1} \left( \sum_{j'=1}^{c} j' w_{j'} + (c+1) \sum_{j'=c+1}^{N} w_{j'} \right).$$

This yields

$$\sum_{k=1}^{N} W_k(X_k, \mathbb{T}_X^{(k)}) \le \sum_{j=1}^{c} w_j + C_d(\delta) \left( \frac{1}{c+1} \sum_{j'=1}^{c} j' w_{j'} + \sum_{j'=c+1}^{N} w_{j'} \right) \le C_d(\delta) + 1.$$

We therefore have

$$\mathbb{E}\left( \sum_{k=1}^{N} W_k(X)^2 \mathbb{E}(Z_k^2 \mid X_k) \right) \le w_1 (C_d(\delta) + 1) \mathbb{E}(h(X)) \to 0,$$

which concludes the proof.                                                                            ∎

Theorem 9.1 is proved in Stone [183] with weaker hypotheses allowing for more flexibility in the computation of distances, in which, for example, differences $X - X_i$ can be normalized by dividing them by a factor $\sigma_i$ that may depend on the training set. These relaxed assumptions slightly complicate the proof, and we refer the reader to Stone [183] for a complete exposition.

### 9.1.2 Optimality

The NN method can be shown to be optimal over some classes of functions. Optimality is in the min-max sense, and works as follows. We assume that the regression function $f(x) = \mathbb{E}(Y \mid X = x)$ belongs to some set $\mathcal{F}$ of real-valued functions on $\mathbb{R}^d$. Most of the time, the estimation methods must be adapted to a given choice of $\mathcal{F}$, and various choices have arisen in the literature: classes of functions with $r$ bounded derivatives, Sobolev or related spaces, functions whose Fourier transforms has given properties, etc.

Consider now an estimator of $f$, denoted $\hat{f}_N$, based on a training set of size $N$. We can measure the error by, say:

$$\left\|\hat{f}_N - f\right\|_2 = \left(\int_{\mathbb{R}^d} (\hat{f}_N(x) - f(x))^2 dx\right)^{1/2}$$

Since $\hat{f}_N$ is computed from a random sample, this error is a random variable. One can study, when $b_N \to 0$, the probability

$$\mathbb{P}_f\left(\|\hat{f}_N - f\|_2^2 \geq cb_N\right)$$

for some constant $c$ and, for example, for the model: $Y = f(X) + \text{noise}$. Here, the notation $\mathbb{P}_f$ refers to the model assumption indicating the unobserved function $f$.

The min-max method considers the worst case and computes

$$M_N(c) = \sup_{f \in \mathcal{F}} \mathbb{P}_f\left(\|\hat{f}_N - f\|_2^2 \geq cb_N\right).$$

This quantity now only depends on the estimation algorithm. One defines the notion of "lower convergence rate" as a sequence $b_N$ such that, for any choice of the estimation algorithm, $M_N(c)$ can be found arbitrarily close to 1 (i.e., $\|\hat{f}_N - f\|_2^2 \geq cb_N$ with arbitrarily high probability for all $f \in \mathcal{F}$), for arbitrarily large $N$ (and for some choice of $c$). The mathematical statement is

$$\exists c > 0 : \liminf_{N \to \infty} M_N(c) = 1.$$

So, if $b_N$ is a lower convergence rate, then, for every estimator, there exists a constant $c$ such that the accuracy $cb_N$ cannot be achieved.

On the other hand, one says that $b_N$ is an achievable rate of convergence if there exists an estimator such that, for some $c'$,

$$\limsup_{N \to \infty} M_N(c') = 0.$$

This says that for large $N$, and for some $c'$, the accuracy is higher than $c'b_N$ for the given estimator. Notice the difference: a lower rate holds for all estimators, and an achievable rate for at least one estimator.

The final definition of a min-max optimal rate is that it is both a lower rate and an achievable rate (obviously for different constants $c$ and $c'$). And an estimator is optimal in the min-max sense if it achieves an optimal rate.

One can show that the $p$-NN estimator is optimal (under some assumptions on the ratio $p_N/N$) when $\mathcal{F}$ is the class of Lipschitz functions on $\mathbb{R}^d$, i.e., the class of functions such that there exists a constant $K$ with

$$|f(x) - f(y)| \leq K|x - y|$$

for all $x, y \in \mathbb{R}^d$. In this case, the optimal rate is $b_N = N^{-1/(2+d)}$ (notice again the "curse of dimensionality": to achieve a given accuracy in the worst case, the number of data points must grow exponentially with the dimension).

If the function class consists of smoother functions (for example, several derivatives), the $p$-NN method is not optimal. This is because the local averaging method is too crude when one knows already that the function is smooth. But it can be modified (for example by fitting, using least squares, a polynomial of some degree instead of computing an average) in order to obtain an optimal rate.

## 9.2   $p$-NN classification

Let $(x_1, y_1, \ldots, x_N, y_N)$ be the training set, with $x_i \in \mathbb{R}^d$ and $y_i \in \mathcal{R}_Y$ where $\mathcal{R}_Y$ is a finite set of classes. Using the same notation as in the previous section, define

$$\widehat{\pi}_w(y|x) = \sum_{k=1}^{N} W_k(x)\mathbf{1}_{y_k=y}.$$

Let the corresponding classifier be

$$\hat{f}_w(x) = \operatorname*{argmax}_{y \in \mathcal{R}_Y} \widehat{\pi}_w(y|x).$$

Theorem 9.1 may be applied, for $y \in \mathcal{R}_Y$, to the function $f_y(x) = \pi(y \mid x) = \mathbb{E}(\mathbf{1}_{Y=y} \mid X = x)$, which allows one to interpret the estimator $\widehat{\pi}(y \mid x)$ as a nearest-neighbor predictor of the random variable $\mathbf{1}_{Y=y}$ as a function of $X$. We therefore obtain the consistency of the estimated posteriors when $N \to \infty$ under the same assumption as those of theorem 9.1. This implies that, for large $N$, the classification will be close to Bayes's rule.

An asymptotic comparison with Bayes's rule can already be made with $p = 1$. Let $\hat{y}_N(x)$ be the 1-NN estimator of $Y$ given $x$ and a training set of size $N$, and let $\hat{y}(x)$ be the Bayes estimator. We can compute the Bayes error by

$$
\begin{aligned}
\mathbb{P}(\hat{y}(X) \neq Y) &= 1 - \mathbb{P}(\hat{y}(X) = Y) \\
&= 1 - \mathbb{E}(\mathbb{P}(\hat{y}(X) = Y | X)) \\
&= 1 - \mathbb{E}(\max_{y \in \mathcal{R}_Y} \pi(y | X))
\end{aligned}
$$

For the 1-NN rule, we have

$$
\begin{aligned}
\mathbb{P}(\hat{y}_N(X) \neq Y) &= 1 - \mathbb{P}(\hat{y}_N(X) = Y) \\
&= 1 - \mathbb{E}(\mathbb{P}(\hat{y}_N(X) = Y | X))
\end{aligned}
$$

Let us make the assumption that nearest neighbors are not tied (with probability one). Let $k^*(x, T)$ denote the index of the nearest neighbor to $x$ in the training set $T$. We have

$$
\begin{aligned}
\mathbb{P}(\hat{y}_N(X) = Y \mid X) &= \mathbb{E}(\mathbb{P}(\hat{y}_N(X) = Y \mid X, \mathbb{T})) \\
&= \mathbb{E}\left( \sum_{k=1}^{N} \mathbb{P}(Y = Y_k \mid X, \mathbb{T}) \mathbf{1}_{k^*(X,\mathbb{T})=k} \right) \\
&= \mathbb{E}\left( \sum_{k=1}^{N} \mathbb{P}(Y = Y_k \mid X, X_k) \mathbf{1}_{k^*(X,\mathbb{T})=k} \right) \\
&= \mathbb{E}\left( \sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} \mathbb{P}(Y = g, Y_k = g \mid X, X_k) \mathbf{1}_{k^*(X,\mathbb{T})=k} \right) \\
&= \mathbb{E}\left( \sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} \pi(g \mid X) \pi(g \mid X_k) \chi_{k^*(X,\mathbb{T})=k} \right) \\
&= \mathbb{E}\left( \sum_{g \in \mathcal{R}_Y} \pi(g \mid X) \pi(g \mid X_{k^*(X,\mathbb{T})}) \right)
\end{aligned}
$$

Now, assume the continuity of $x \mapsto \pi(g \mid x)$ (although the result can be proved without this simplifying assumption). We know that $X_{k^*(X,\mathbb{T})} \to X$ when $N \to \infty$ (see the proof of theorem 9.1), which implies that $\pi(g \mid X_{k^*(X,\mathbb{T})}) \to \pi(x \mid X)$ and at the limit

$$
\mathbb{P}(\hat{y}_N(X) = Y \mid X) \to \sum_{g \in \mathcal{R}_Y} \pi(g \mid X)^2.
$$

This implies that the asymptotic 1-NN misclassification error is always smaller than 2 times the Bayes error, that is

$$1 - \mathbb{E}\left( \sum_{g \in \mathcal{R}_Y} \pi(g \mid X)^2 \right) \leq 2(1 - \mathbb{E}(\max_g \pi(g \mid X)))$$

Indeed, the left-hand term is smaller than $1 - \mathbb{E}(\max_g \pi(g|x)^2)$ and the result comes from the fact that for any $t \in \mathbb{R}$. $1 - t^2 \leq 2 - 2t$.

**Remark 9.2** Nearest neighbor methods may require large computation time, since, for a given $x$, the number of comparisons which are needed is the size of the training set. However, efficient (tree-based) search algorithms can be used in many cases to reduce it to a logarithm in the size of the database, which is acceptable. A reduction of the size of the training set by clustering also is a possibility for improving the efficiency.

The computation time is also generally proportional to the dimension $d$ of the input $x$. When $d$ is large, a reduction of dimension is often a good idea. Principal components (see chapter 20), or LDA directions (see chapter 8) can be used for this purpose. ♦

## 9.3   Designing the distance

**LDA-based distance**   The most important factor in the design of a NN procedure probably is the choice of the distance, something we have not discussed so far. Intuitively, the distance should increase fast in the directions "perpendicular" to the regions of constancy of the class variables, and slowly (ideally not at all) within these regions. The following construction uses discriminant analysis [87].

For $g \in \mathcal{R}_Y$, let $\Sigma_g$ be the covariance matrix in class $g$, and $\Sigma_w = \sum_{g \in \mathcal{R}_Y} \pi_g \Sigma_g$ be the within-class variance, where $\pi_g$ is the frequency of class $g$. Let $\Sigma_b$ denote the between-class covariance matrix (see section 8.2).

For $x \in \mathbb{R}^d$, define the spherized vector $x^* = \Sigma_w^{-1/2} x$. The between-class variance computed for spherized data is $\Sigma_b^* = \Sigma_w^{-1/2} \Sigma_b \Sigma_w^{-1/2}$. A direction is discriminant if it is close to the principal eigenvectors of $\Sigma_b^*$. This suggests the introduction of the norm

$$|x|_*^2 = (x^*)^T \Sigma_b^* x^* = x^T \Sigma_w^{-1/2} (\Sigma_w^{-1/2} \Sigma_b \Sigma_w^{-1/2}) \Sigma_w^{-1/2} x = x^T \Sigma_w^{-1} \Sigma_b \Sigma_w^{-1} x.$$

This replaces the standard Euclidean norm (the method can be made more robust by adding $\epsilon \mathrm{Id}_{\mathbb{R}^d}$ to $\Sigma_b^*$.)

**Tangent distance** Designing the distance, however, can sometimes be based on *a priori* knowledge on some invariance properties associated with the classes. A successful example comes from character recognition, where it is known that transforming images by slightly rotating, scaling, or translating the character should not change its class. This corresponds to the following general framework.

For each input $x \in \mathbb{R}^d$, assume that one can make small transformations without changing the class of $x$. We model these transformations as parametrized functions $x \mapsto x_\theta = \varphi(x, \theta) \in \mathbb{R}^d$, such that $\varphi(x, 0) = x$ and $\varphi$ is smooth in $\theta$, which is a $q$-dimensional parameter. The assumption is that $\varphi(x, \theta)$ and $x$ should be from the same class, at least for small $\theta$. This will be used to improve on the Euclidean distance on $\mathbb{R}^d$.

Take $x, x' \in \mathbb{R}^d$. Ideally, one would like to use the distance $D(x, x') = \inf_{\theta, \theta'} \operatorname{dist}(x_\theta, x_{\theta'})$ where $\theta$ and $\theta'$ are restricted to a small neighborhood of 0. A more tractable expression can be based on first-order approximations

$$x_\theta \simeq x + \nabla_\theta \varphi(x, 0)u = x + \sum_{i=1}^{q} u_i \partial_{\theta_i} \varphi(x, 0)$$

$$\text{and} \quad x'_\theta \simeq x' + \nabla_\theta \varphi(x', 0)u' = x' + \sum_{i=1}^{q} u'_i \partial_{\theta_i} \varphi(x', 0)$$

yielding the approximation (also called the tangent distance)

$$D(x, x')^2 \simeq \inf_{u, u' \in \mathbb{R}^q} \left\| x - x' + \nabla_\theta \varphi(x, 0)u - \nabla_\theta \varphi(x', 0)u' \right\|^2.$$

The computation now is a simple least-squares problem, for which the solution is given by the system

$$\begin{pmatrix} \nabla_\theta \varphi(x, 0)^T \nabla_\theta \varphi(x, 0) & -\nabla_\theta \varphi(x, 0)^T \nabla_\theta \varphi(x', 0) \\ -\nabla_\theta \varphi(x', 0)^T \nabla_\theta \varphi(x, 0) & \nabla_\theta \varphi(x', 0)^T \nabla_\theta \varphi(x', 0) \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \nabla_\theta \varphi(x, 0)^T (x' - x) \\ \nabla_\theta \varphi(x', 0)^T (x - x') \end{pmatrix}.$$

A slight modification, to ensure that the norms of $u$ and $u'$ are not too large, is to add a penalty $\lambda(|u|^2 + |u'|^2)$, which results in adding $\lambda \operatorname{Id}_{\mathbb{R}^q}$ to the diagonal blocs of the above matrix.

# Chapter 10

# Tree-based Algorithms, Randomization and Boosting

## 10.1 Recursive Partitioning

Recursive partitioning methods implement a "divide and conquer" strategy to address the prediction problem. They separate the input space $\mathcal{R}_X$ into small regions on which prediction is "easy," i.e., such that the observed values of the output variable are (almost) constant for input values in these regions. The regions are estimated by recursive divisions until they become either too small or homogeneous. These divisions are conveniently represented in the form of binary trees.

### 10.1.1 Binary prediction trees

Define a *binary node* to be a structure $v$ that contains the following information (note that the definition is recursive):

- A *label* $L(v)$ that uniquely identifies the node.

- A *set of children*, $C(v)$, which is either empty or a pair of nodes $(l(v), r(v))$.

- A *binary feature*, i.e., a function $\gamma_v : \mathcal{R}_X \rightarrow \{0, 1\}$, which is "None" (i.e., irrelevant) if the node has no children.

- A *predictor*, $f_v : \mathcal{R}_X \rightarrow \mathcal{R}_Y$, which is "None" if the node has children.

A node without children is called a *terminal node*, or a leaf.

A *binary prediction tree* $\mathfrak{T}$ is a finite set of nodes, with the following properties:

(i) Only one node has no parent (the root, denoted $\rho$ or $\rho_{\mathfrak{T}}$);

(ii)  Each other node has exactly one parent;

(iii)  No node is a descendent of itself.

### 10.1.2   Training algorithm

Assume that a family $\Gamma$ of binary features $\gamma : \mathcal{R}_X \to \{0,1\}$ is chosen, together with a family $\mathcal{F}$ of predictors $f : \mathcal{R}_X \to \mathcal{R}_Y$. Assume also the existence of two "algorithms" as follows:

• Feature selection: Given the feature set $\Gamma$ and a training set $T$, return an optimized binary feature $\widehat{\gamma}_{T,\Gamma} \in \Gamma$.

• Predictor optimization: Given the predictor set $\mathcal{F}$ and a training set $T$, return an optimized predictor $\hat{f}_{T,\mathcal{F}} \in \mathcal{F}$.

Finally, assume that a stopping rule is defined, as a function of training sets $\sigma : T \mapsto \sigma(T) \in \{0,1\}$, where 0 means "continue", and 1 means "stop".

Given a training set $T_0$, the algorithm builds a binary tree $\mathfrak{T}$ using a recursive construction. Each node $v \in \mathfrak{T}$ will be associated to a subset of $T_0$, denoted $T_v$. We define below a recursive operation, denoted Node$(T, j)$ that adds a node $v$ to a tree $\mathfrak{T}$ given a subset $T$ of $T_0$ and a label $j$. Starting with $\mathfrak{T} = \emptyset$, calling Node$(T_0, 0)$ will then create the desired tree.

---

**Algorithm 10.1 (Node insertion: Node$(T, j)$)**

(a)  Given $T$ and $j$, let $T_v = T$ and $L(v) = j$.

(b)  If $\sigma(T) = 1$, let $C(v) = \emptyset$, $\gamma_v =$ "None" and $f_v = \hat{f}_{T,\mathcal{F}}$.

(c)  If $\sigma(T) = 0$, let $f_v =$ "None", $\gamma_v = \widehat{\gamma}_{T,\Gamma}$ and $C(v) = (l(v), r(v))$ with

$$l(v) = \text{Node}(T_l, 2j + 1), \quad r(v) = \text{Node}(T_r, 2j + 2)$$

where
$$T_l = \{(x,y) \in T : \gamma_v(x) = 0\}, \quad T_r = \{(x,y) \in T : \gamma_v(x) = 1\}$$

(d)  Add $v$ to $\mathfrak{T}$ and return.

---

**Remark 10.1**  Note that, even though the learning algorithm for prediction trees can be very conveniently described in recursive form as above, efficient computer implementations should avoid recursive calls, which may be inefficient and memory demanding. Moreover, for large trees, it is likely that recursive implementations will reach the maximal number of recursive calls imposed by compilers.                    ♦

### 10.1.3 Resulting predictor

Once the tree is built, the predictor $x \mapsto \hat{f}_{\mathbb{T}}(x)$ is recursively defined as follows.

(a) Initialize the computation with $v = \rho$.

(b) At a given step of the algorithm, let $v$ be the current node.

- If $v$ has no children: then let $\hat{f}_{\mathbb{T}}(x) = f_v(x)$.
- Otherwise: replace $v$ by $l(v)$ if $\gamma_v(x) = 0$ and by $r(v)$ if $\gamma_v(x) = 1$ and go back to (b).

### 10.1.4 Stopping rule

The function $\sigma$, which decides whether a node is terminal or not is generally defined based on very simple rules. Typically, $\sigma(T) = 1$ when one the following conditions is satisfied:

- The number of training examples in $T$ is small (e.g., less than 5).

- The values $y_k$ in $T$ have a small variance (regression) or are constant (classification).

### 10.1.5 Leaf predictors

When one reaches a terminal node $v$ (so that $\sigma(T_v) = 1$), a predictor $f_v$ must be determined. This function can be optimized within any set $\mathcal{F}$ of predictors, using any learning algorithm, but in practice, one usually makes this fairly simple and defines $\mathcal{F}$ to be the family of constant functions taking values in $\mathcal{R}_Y$. The function $\hat{f}_{T,\mathcal{F}}$ is then defined as:

- the average of the values of $y_k$, for $(x_k, y_k) \in T$ (regression);

- the mode of the distribution of $y_k$, for $(x_k, y_k) \in T$ (classification).

### 10.1.6 Binary features

The space $\Gamma$ of possible binary features must be specified in order to partition non-terminal nodes. A standard choice, used in the CART model [43] with $\mathcal{R}_X = \mathbb{R}^d$, is

$$\Gamma = \left\{ \gamma(x) = \mathbf{1}_{[x^{(i)} \geq \theta]}, i = 1, \ldots, d, \theta \in \mathbb{R} \right\} \tag{10.1}$$

where $x^{(i)}$ is the $i$th coordinate of $x$. This corresponds to splitting the space using a hyperplane parallel to one of the coordinate axes.

The binary function $\widehat{\gamma}_{T,\Gamma}$ can be optimized over $\Gamma$ using a greedy evaluation of the risk, assuming that the prediction is based on the two nodes resulting from the split. For $\gamma \in \Gamma$, $f_0, f_1 \in \mathcal{F}$, define

$$F_{\gamma, f_0, f_1}(x) = \begin{cases} f_0(x) \text{ if } \gamma(x) = 0 \\ f_1(x) \text{ if } \gamma(x) = 1 \end{cases}$$

Given a risk function $r$, one then evaluates

$$\mathcal{E}_T(\gamma) = \min_{f_0, f_1 \in \mathcal{F}} \sum_{(x,y) \in T} r(y, F_{\gamma, f_0, f_1}(x))$$

One then chooses $\widehat{\gamma}_{T,\Gamma} = \operatorname{argmin}_{\gamma \in \Gamma}(\mathcal{E}_T(\Gamma))$.

**Example 10.2 (Regression)** Consider the regression case, taking squared differences as risk and letting $\mathcal{F}$ contain only constant functions. Then

$$\mathcal{E}_T(\gamma) = \min_{m_0, m_1} \sum_{(x,y) \in T} \left( (y - m_0)^2 \mathbf{1}_{\gamma(x)=0} + (y - m_1)^2 \mathbf{1}_{\gamma(x)=1} \right).$$

Obviously, the optimal $m_0$ and $m_1$ are the averages of the output values, $y$, in each of the subdomains defined by $\gamma$. For CART (see (10.1)), this cost must be minimized over all choices $(i, \theta)$ with $i = 1, \dots, d$ and $\theta \in \mathbb{R}$ where $\gamma_{i,\theta}(x) = 1$ if $x(i) > \theta$ and 0 otherwise.                                                                          ◆

**Example 10.3 (Classification.)** For classification, one can apply the same method, with the 0/1 loss, letting

$$\mathcal{E}_T(\gamma) = \min_{g_0, g_1} \sum_{(x,y) \in T} \left( \mathbf{1}_{y \neq g_0} \mathbf{1}_{\gamma(x)=0} + \mathbf{1}_{y \neq g_1} \mathbf{1}_{\gamma(x)=1} \right).$$

The optimal $g_0$ and $g_1$ are the majority classes in $T \cap \{\gamma = 0\}$ and $T \cap \{\gamma = 1\}$.         ◆

**Example 10.4 (Entropy selection for classification)** For classification trees, other splitting criteria may be used based on the empirical probability $p_T$ on the set $T$, defined as

$$p_T(A) = \frac{1}{N} |\{k : (x_k, y_k) \in A\}|$$

for $A \subset \mathcal{R}_X \times \mathcal{R}_Y$. The previous criterion, $\mathcal{E}_T(\gamma)$, is proportional to

$$p_T(\gamma = 0)(1 - \max_g p_T(g \mid \gamma = 0)) + p_T(\gamma = 1)(1 - \max_g p_T(g \mid \gamma = 1)).$$

One can define alternative objectives in the form

$$p_T(\gamma = 0)\mathcal{H}(p_T(g \mid \gamma = 0)) + p_T(\gamma = 1)\mathcal{H}(p_T(g \mid \gamma = 1))$$

where $\pi \to \mathcal{H}(\pi)$ associates to a probability distribution $\pi$ a "complexity measure" that is minimal when $\pi$ is concentrated on a single class (which is the case for $\pi \mapsto 1 - \max_g \pi(g)$).

Many such measures exists, and many of them are defined as various forms of entropy designed in information theory. The most celebrated is Shannon's entropy [176], defined by

$$\mathcal{H}(p) = -\sum_{g \in \mathcal{R}_Y} p(g) \log p(g).$$

It is always positive, and minimal when the distribution is concentrated on a single class. Other entropy measures include:

- The Tsallis entropy: $\mathcal{H}(p) = \frac{1}{1-q} \sum_{g \in \mathcal{R}_Y} (p(g)^q - 1)$, for $q \neq 1$. (Tsallis entropy for $q = 2$ is sometimes called the Gini impurity index.)
- The Renyi entropy: $\mathcal{H}(p) = \frac{1}{1-q} \log \sum_{g \in \mathcal{R}_Y} p(g)^q$, for $q \geq 0, q \neq 1$. ◆

### 10.1.7 Pruning

Growing a decision tree to its maximal depth (given the amount of available data) generally leads to predictors that overfit the data. The training algorithm is usually followed by a pruning step that removes some some nodes based on a complexity penalty.

Letting $\tau(\mathfrak{T})$ denote the set of terminal nodes in the tree $\mathfrak{T}$ and $\hat{f}_{\mathfrak{T}}$ the associated predictor, pruning is represented as an optimization problem, where one minimizes, given the training set $T$,

$$U_\lambda(\mathfrak{T}, T) = \hat{R}_T(\hat{f}_{\mathfrak{T}}) + \lambda |\tau(\mathfrak{T})|$$

where $\hat{R}_T$ is as usual the in-sample error measured on the training set $T$.

To prune a tree, one selects one or more internal nodes and remove all their descendants (so that these nodes become terminal). Associate to each node $v$ in $\mathfrak{T}$ its local in-sample error $\mathcal{E}_{T_v}$ equal to the error made by the optimal classifier estimated from the training data associated with $v$. Then,

$$U_\lambda(\mathfrak{T}, T) = \sum_{v \in \tau(\mathfrak{T})} \frac{|T_v|}{|T|} \mathcal{E}_{T_v} + \lambda |\tau(\mathfrak{T})|$$

If $v$ is a node in $\mathfrak{T}$ (internal or terminal), let $\mathfrak{T}_v$ be the subtree of $\mathfrak{T}$ containing $v$ as a root and all its descendants. Let $\mathfrak{T}^{(v)}$ be the tree $\mathfrak{T}$ will all descendants of $v$

removed (keeping $v$). Then

$$U_\lambda(\mathfrak{T}, T) = U_0(\mathfrak{T}^{(v)}, T) - \frac{|T_v|}{|T|}(\mathcal{E}_{T_v} - U_0(\mathfrak{T}_v, T_v)) + \lambda(|\tau(\mathfrak{T}_v)| - 1).$$

Note also that, if $v$ is internal, and $v'$, $v''$ are its children, then

$$U_0(\mathfrak{T}_v, T_v) = \frac{|T_{v'}|}{|T_v|}U_0(\mathfrak{T}_{v'}, T_{v'}) + \frac{|T_{v''}|}{|T_v|}U_0(\mathfrak{T}_{v''}, T_{v''})$$

This formula can be used to compute $U_0(\mathfrak{T}_v)$ recursively for all nodes, starting with leaves for which $U_0(\mathfrak{T}_v) = \mathcal{E}(T_v)$. (We also have $|\tau(\mathfrak{T}_v)| = |\tau(\mathfrak{T}_{v'})| + |\tau(\mathfrak{T}_{v''})|$.) The following algorithm converges to a *global minimizer* of $U_\lambda$.

---

**Algorithm 10.2 (Pruning)**

(1)  Start with a complete tree $\mathfrak{T}(0)$ built without penalty.

(2)  Compute, for all nodes $U_0(\mathfrak{T}_v)$ and $|\tau(\mathfrak{T}_v)|$. Let

$$\psi_v = \frac{|T_v|}{|T|}(\mathcal{E}_{T_v} - U_0(\mathfrak{T}_v)) - \lambda(|\tau(\mathfrak{T}_v)| - 1).$$

(3)  Iterate the following steps.

- If $\psi_v < 0$ for all internal nodes $v$, exit the program and return the current $\mathfrak{T}(n)$.
- Otherwise choose an internal node $v$ such that $\psi_v$ is largest.
- Let $\mathfrak{T}(n+1) = \mathfrak{T}^{(v)}(n)$. Subtract $\lambda(|\tau(\mathfrak{T}_v(n))| - 1)$ to $\rho_{v'}$ for all $v'$ ancestor of $v$.

---

## 10.2   Random Forests

### 10.2.1   Bagging

A random forest [7, 42] is a special case of composite predictors (we will see other examples later in this chapter when describing boosting methods) that train multiple individual predictors under various conditions and combine them, through averaging, or majority voting. With random forests, one generates individual trees by randomizing the parameters of the learning process. One way to achieve this is to randomly sample from the training set before running the training algorithm.

Letting as before $T_0 = (x_1, y_1, \ldots, x_N, y_N)$ denote the original set, with size $N$, one can create "new" training data by sampling with replacement from $T_0$. More precisely, consider the family of independent random variables $\xi = (\xi_1, \ldots, \xi_N)$, with

each $\xi_j$ following a uniform distribution over $\{1,\dots,N\}$. One can then form the random training set

$$T_0(\xi) = (x_{\xi_1}, y_{\xi_1}, \dots, x_{\xi_N}, y_{\xi_N}).$$

Running the training algorithm using $T_0(\xi)$ then provides a random tree, denoted $\mathbb{T}(\xi)$. Now, by sampling $K$ realizations of $\xi$, say $\xi^{(1)},\dots,\xi^{(K)}$, one obtains a collection of $K$ random trees (a random forest) $\mathbb{T}^* = (\mathbb{T}_1,\dots,\mathbb{T}_K)$, with $\mathbb{T}_j = \mathbb{T}(\xi^{(j)})$ that can be combined to provide a final predictor. The simplest way to combine them is to average the predictors returned by each tree (assuming, for classification, that this predictor is a probability distribution on classes), so that

$$f_{\mathbb{T}^*}(x) = \frac{1}{K} \sum_{j=1}^{K} f_{\mathbb{T}_j}(x). \tag{10.2}$$

For classification, one can alternatively let each individual tree "vote" for their most likely class.

Obviously, randomizing training data and averaging the predictors is a general approach that can be applied to any prediction algorithm, not only to decision trees. In the literature, the approach described above has been called *bagging* [41], which is an acronym for "bootstrap aggregating" (bootstrap itself being a general resampling method in statistics that samples training data with replacement to determine some properties of estimators). Another way to randomize predictors (especially when $d$, the input dimension is large), is to randomize input data by randomly removing some of the coordinates, leading to a similar construction.

With decision trees one can in addition randomize the binary features use to split nodes, as described next. While bagging may provide some enhancement to predictors, feature randomization for decision trees often significantly improves the performance, and is the typical randomization method used for random forests.

### 10.2.2  Feature randomization

When one decides to split a node during the construction of a prediction tree, one can optimize the binary feature $\gamma$ over a random subset of $\Gamma$ rather than exploring the whole set. For CART, for example, one can select a small number of dimensions $i_1,\dots,i_q \in \{1,\dots,d\}$ with $q \ll d$, and optimize $\gamma$ by thresholding one of the coordinates $x^{(i_j)}$ for $j \in \{1,\dots,q\}$. This results in a randomized version of the node insertion function.

---

**Algorithm 10.3 (Randomized node insertion: RNode$(T,j)$)**

   (a) Given $T$ and $j$, let $T_v = T$ and $L(v) = j$.

   (b) If $\sigma(T) = 1$, let $C(v) = \emptyset$, $\gamma_v = $ "None" and $f_v = \hat{f}_{T,CF}$.

(c) If $\sigma(T') = 0$, *sample (e.g., uniformly without replacement) a subset* $\Gamma_v$ *of* $\Gamma$ *and let* $f_v = $ "None", $\gamma_v = \hat{\gamma}_{T,\Gamma_v}$ and $C(v) = (l(v), r(v))$ with

$$l(v) = \mathrm{Node}(T_l, 2j+1)$$
$$r(v) = \mathrm{Node}(T_r, 2j+2)$$

where

$$T_l = \{(x, y) \in T : \gamma_v(x) = 0\}$$
$$T_r = \{(x, y) \in T : \gamma_v(x) = 1\}$$

(d) Add $v$ to $\mathfrak{T}$ and return.

---

Now, each time the function $\mathrm{RNode}(T_0, 0)$ is run, it returns a different, random, tree. If it is called $K$ times, this results in a random forest $\mathfrak{T}^* = (\mathfrak{T}_1, \dots \mathfrak{T}_K)$, with a predictor $\mathcal{F}_{\mathfrak{T}^*}$ given by (10.2). Note that trees in random forests are generally not pruned, since this operation has been observed to bring no improvement in the context of randomized tress.

## 10.3   Top-Scoring Pairs

Top-Scoring Pair (TSP) classifiers were introduced in Geman et al. [78] and can be seen as forests formed with depth-one classification trees in which splitting rules are based on the comparison of pairs of variables. More precisely, define

$$\gamma_{ij}(x) = \mathbf{1}_{x^{(i)} > x^{(j)}}.$$

A decision tree based on these rules only relies on the order between the features, and is therefore well adapted to situations in which the observations are subject to increasing transformations, i.e., when the observed variable $X$ is such that $X^{(j)} = \varphi(Z^{(j)})$, where $\varphi : \mathbb{R} \to \mathbb{R}$ is random and increasing and $Z$ is a latent (unobserved) variable. Obviously, in such a case, order-based splitting rules do not depend on $\varphi$. Such an assumption is relevant, for example, when experimental conditions (such as temperature) may affect the actual data collection, without changing their order, which is the case when measuring high-throughput biological data, such as microarrays, for which the approach was introduced.

Assuming two classes, a depth-one tree in this context is simply the classifier $f_{ij} = \gamma_{ij}$. Given a training set, the associated empirical error is

$$\mathcal{E}_{ij} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{1}_{\gamma_{ij}(x_k) \neq y_k} = \frac{1}{N} \sum_{k=1}^{N} |y_k - \gamma_{ij}(x_k)|$$

and the balanced error (better adapted to situations in which one class is observed more often than the other) is

$$\mathcal{E}_{ij}^b = \sum_{k=1}^{N} w_k |y_k - \gamma_{ij}(x_k)|$$

with $w_k = 1/(2N_{y_k})$, where $N_0$, $N_1$ are the number of observations with $y_k = 0$, $y_k = 1$. Pairs $(i, j)$ with small errors are those for which the order between the features switch with high probability when passing from class 0 to class 1.

In its simplest form, the TSP classifier defines the set

$$\mathcal{P} = \underset{ij}{\operatorname{argmin}} \, \mathcal{E}_{ij}^b$$

of global minimizers of the empirical error (which may just be a singleton) and predicts the class based on a majority vote among the family of predictors $(f_{ij}, (i, j) \in \mathcal{P})$. Equivalently, selected variables maximize the score $\Delta_{ij} = 1 - \mathcal{E}_{ij}^b$, leading to the method's name.

Such classifiers, which are remarkably simple, have been found to be competitive among a wide range of "advanced" classification algorithms for large-dimensional problems in computational biology. The method has been refined in Tan et al. [190], leading to the $k$-TSP classifier, which addresses the following remarks. First, when $j, j'$ are highly correlated, and $(i, j)$ is a high-scoring pair, then $(i, j')$ is likely to be one too, and their associated decision rules will be redundant. Such cases should preferably be pruned from the classification rules, especially if one wants to select a small number of pairs. Second, among pairs of features that switch with the same probability, it is natural to prefer those for which the magnitude of the switch is largest, e.g., when the pair of variables switches from a regime in which one of them is very low and the other very high to the opposite. In Tan et al. [190], a rank-based tie-breaker is introduced, defined as

$$\rho_{ij} = \sum_{k=1}^{N} w_k (R_k(i) - R_k(j))(2y_k - 1),$$

where $R_k(i)$ denotes the rank of $x_k^{(i)}$ in $x_k^{(1)}, \ldots, x_k^{(d)}$. One can now order pairs $(i, j)$ and $(i', j')$ by stating that the former scores higher if (i) $\Delta_{ij} > \Delta_{i'j'}$, or (ii) $\Delta_{ij} = \Delta_{i'j'}$ and $\rho_{ij} > \rho_{i'j'}$. The $k$-TSP classifier is formed by selecting pairs, starting from the highest scoring one, and use as $l$th pair (for $l \leq k$) the highest scoring ones among all those that do not overlap with the previously selected ones. In [190], the value of $k$ is optimized using cross-validation.

## 10.4  Adaboost

Boosting methods refer to algorithms in which classifiers are enhanced by recursively making them focus on harder data. We first address the issue of classification, and describe one of the earliest algorithms (Adaboost). We will then interpret it as a greedy gradient descent algorithm, as this interpretation will lead to further extensions.

### 10.4.1  General set-up

We first consider binary classification problems, with $\mathcal{R}_Y = \{-1, 1\}$. We want to design a function $x \mapsto F(x) \in \{-1, 1\}$ on the basis of a training set $T = (x_1, y_1, \ldots, x_N, y_N)$. With the 0-1 loss, minimizing the empirical error is equivalent to maximizing

$$\mathcal{E}_T(F) = \frac{1}{N} \sum_{k=1}^{N} y_k F(x_k).$$

Boosting algorithms build the function $F$ as a linear combination of "base classifiers," $f_1, \ldots, f_M$, taking

$$F(x) = \operatorname{sign}\left( \sum_{j=1}^{M} \alpha_j f_j(x) \right).$$

We assume that each base classifier, $f_j$, takes values in $[-1, 1]$ (the interval).

The sequence of base classifiers is learned by progressively focusing on the hardest examples. We will therefore assume that the training algorithm for base classifiers takes as input the training set $T$ as well a family of positive weights $W = (w_1, \ldots, w_N)$. More precisely, letting

$$p_W(k) = \frac{w_k}{\sum_{k=1}^{N} w_k},$$

the weighted algorithm should implement (explicitly or implicitly) the equivalent of an unweighted algorithm on a simulated training set obtained by sampling with replacement $K \gg N$ elements of $T$ according to $p_W$ (ideally letting $K \to \infty$). Let us take a few examples.

- Weighted LDA: one can use LDA as described in section 8.2 with

$$c_g = \sum_{k:y_k=g} p_W(k), \quad \mu_g = \frac{1}{c_g} \sum_{k:y_k=g} p_W(k)x_k, \quad \mu = \sum_{g \in \mathcal{R}_Y} c_g \mu_g$$

and the covariance matrices:

$$\Sigma_w = \sum_{k=1}^{N} p_W(k)(x_k - \mu_{y_k})(x_k - \mu_{y_k})^T, \quad \Sigma_b = \sum_{g \in \mathcal{R}_Y} c_g(\mu_g - \bar{\mu})(\mu_g - \bar{\mu})^T.$$

- Weighted logistic regression: just maximize

$$\sum_{k=1}^{N} p_W(k) \log \pi_\theta(y_k | x_k)$$

where $\pi_\theta$ is given by the logistic model.

- Empirical risk minimization algorithms can be modified in order to minimize

$$\hat{R}_{T,W}(f) = \sum_{k=1}^{N} w_k r(y_k, f(x_k)).$$

- Of course, any algorithm can be run on a training set resampled using $p_W$.

### 10.4.2   The Adaboost algorithm

Boosting algorithms keep track of a family of weights and modify it after the $j$th classifier $f_j$ is computed, increasing the importance of misclassified examples, before computing the next classifier. The following algorithm, called Adaboost [172, 73], describes one such approach.

---

**Algorithm 10.4 (Adaboost)**
- Start with uniform weights, letting $W(1) = (w_1(1),\dots,w_N(1))$ with $w_k(1) = 1/N$, $k = 1,\dots,N$. Fix a number $\rho \in (0,1]$ and an integer $M > 0$.

- Iterate, for $j = 1,\dots,M$:

   (1) Fit a base classifier $f_j$ using the weights $W(j) = (w_1(j),\dots,w_N(j))$. Let

$$S_w^+(j) = \sum_{k=1}^{N} w_k(j)(2 - |y_k - f_j(x_k)|) \tag{10.3a}$$

$$S_w^-(j) = \sum_{k=1}^{N} w_k(j)|y_k - f_j(x_k)| \tag{10.3b}$$

and define $\alpha_j = \rho \log\left(S_w^+(j)/S_w^-(j)\right)$
   (2) Update the weights by

$$w_k(j+1) = w_k(j)\exp\left(\alpha_j|y_k - f_j(x_k)|/2\right).$$

- Return the classifier:

$$F(x) = \text{sign}\left(\sum_{j=1}^{M} \alpha_j f_j(x)\right).$$

---

If $f_j$ is binary, i.e., $f_j(x) \in \{-1, 1\}$, then $|y_k - f_j(x_k)| = 2\mathbf{1}_{y_k \neq f_j(x_k)}$, so that $S_W^+/2$ is the weighted number of correct classifications and $S_W^-/2$ is the weighted number of incorrect ones.

For $\alpha_j$ to be positive, the $j$th classifier must do better than pure chance on the weighted training set. If not, taking $\alpha_j \leq 0$ reflects the fact that, in that case, $-f_j$ has better performance on training data.

Algorithms that do slightly better than chance with high probability are called "weak learners" [172]. The following proposition [73] shows that, if the base classifiers reliably perform strictly better than chance (by a fixed, but not necessarily large, margin), then the boosting algorithm can make the training-set error arbitrarily close to 0.

**Proposition 10.5** *Let $\mathcal{E}_T$ be the training set error of the estimator $F$ returned by Algorithm 10.4, i.e.,*

$$\mathcal{E}_T = \frac{1}{N} \sum_{k=1}^{N} \mathbf{1}_{y_k \neq F(x_k)}.$$

*Then*

$$\mathcal{E}_T \leq \prod_{j=1}^{M} \left(\epsilon_j^\rho (1 - \epsilon_j)^{1-\rho} + \epsilon_j^{1-\rho}(1 - \epsilon_j)^\rho\right)$$

*where*

$$\epsilon_j = \frac{S_W^-(j)}{S_W^+(j) + S_W^-(j)}.$$

Proof We note that example $k$ is misclassified by the final classifier if and only if

$$\sum_{j=1}^{M} \alpha_j y_k f_j(x_k) \leq 0$$

or

$$\prod_{j=1}^{M} e^{-\alpha_j y_k f_j(x_k)/2} \geq 1$$

Noting that $|y_k - f_j(x_k)| = 1 - y_k f_j(x_k)$, we see that example $k$ is misclassified when

$$\prod_{j=1}^{M} e^{\alpha_j |y_k - f_j(x_k)|/2} \geq \prod_{j=1}^{M} e^{\alpha_j/2}.$$

This shows that

$$
\begin{aligned}
\mathcal{E}_T &= \frac{1}{N} \sum_{k=1}^{N} \mathbf{1}_{y_k \neq F(x_k)} \\
&= \frac{1}{N} \sum_{k=1}^{N} \mathbf{1}_{\prod_{j=1}^{M} e^{\alpha_j |y_k - f_j(x_k)|/2} \geq \prod_{j=1}^{M} e^{\alpha_j/2}} \\
&\leq \frac{1}{N} \sum_{k=1}^{N} \prod_{j=1}^{M} e^{\alpha_j |y_k - f_j(x_k)|/2} \prod_{j=1}^{M} e^{-\alpha_j/2}.
\end{aligned}
$$

Let, for $q \leq M$,

$$U_q = \frac{1}{N} \sum_{k=1}^{N} \prod_{j=1}^{q} e^{\alpha_j |y_k - f_j(x_k)|/2}.$$

Since

$$w_k(q) = \frac{1}{N} \prod_{j=1}^{q-1} e^{\alpha_j |y_k - f_j(x_k)|/2},$$

we also have $U_q = \sum_{k=1}^{N} w_k(q+1) = (S_W^+(q+1) + S_W^-(q+1))/2$.

We will use the inequality [1]

$$e^{\alpha t} \leq 1 - (1 - e^{\alpha})t,$$

---

[1]This inequality is clear for $\alpha = 0$. Assuming $\alpha \neq 0$, the difference between the upper and lower bound is

$$q(t) = 1 - e^{\alpha t} - (1 - e^{\alpha})t.$$

The function $q$ is concave (its second derivative is $-\alpha^2 e^{\alpha t}$) with $q(0) = q(1) = 0$ and therefore non-negative over $[0, 1]$.

which is true for all $\alpha \in \mathbb{R}$ and $t \in [0,1]$, to write

$$
\begin{aligned}
U_q &\leq \frac{1}{N} \sum_{k=1}^{N} \prod_{j=1}^{q-1} e^{\alpha_j |y_k - f_j(x_k)|/2} (1 - (1 - e^{\alpha_q})|y_k - f_q(x_k)|/2) \\
&= \sum_{k=1}^{N} w_k(q)(1 - (1 - e^{\alpha_q})|y_k - f_q(x_k)|/2) \\
&= \sum_{k=1}^{N} w_k(q) - (1 - e^{\alpha_q}) \sum_{k=1}^{N} w_k(q)|y_k - f_q(x_k)|/2 \\
&= U_{q-1}(1 - (1 - e^{\alpha_q})\epsilon_q)
\end{aligned}
$$

This gives (using $U_0 = 1$)

$$
U_M \leq \prod_{j=1}^{M} \left( 1 - \left(1 - e^{\alpha_j}\right)\epsilon_j \right)
$$

and

$$
\mathcal{E}_T \leq \prod_{j=1}^{M} \left( 1 - \left(1 - e^{\alpha_j}\right)\epsilon_j \right) e^{-\alpha_j/2}.
$$

It now suffices to replace $e^{\alpha_j}$ by $(1 - \epsilon_j)^\rho \epsilon_j^{-\rho}$ and note that

$$
\left( 1 - (1 - (1 - \epsilon_j)^\rho \epsilon_j^{-\rho})\epsilon_j \right)(1 - \epsilon_j)^{-\rho/2} \epsilon_j^{\rho/2} = \epsilon_j^\rho (1 - \epsilon_j)^{1-\rho} + \epsilon_j^{1-\rho}(1 - \epsilon_j)^\rho
$$

to conclude the proof.                                                                     ∎

For $\epsilon \in [0,1]$, one has

$$
\epsilon^\rho (1 - \epsilon)^{1-\rho} + \epsilon^{1-\rho}(1 - \epsilon)^\rho = 1 - (\epsilon^\rho - (1 - \epsilon)^\rho)(\epsilon^{1-\rho} - (1 - \epsilon)^{-1-\rho}) \leq 1
$$

with equality if and only if $\epsilon = 1/2$, so that each term in the upper-bound reduces the error unless the corresponding base classifier does not perform better than pure chance. The parameter $\rho$ determines the level at which one increases the importance of misclassified examples for the next step. Let $\tilde{S}_W^+(j)$ and $\tilde{S}_W^-(j)$ denote the expressions in (10.3a) and (10.3b) with $w_k(j)$ replaced by $w_k(j+1)$. Then, in the case when the base classifiers are binary, ensuring that $|y_k - f_j(x_k)|/2 = \mathbf{1}_{y_k \neq f_j(x_k)}$, one can easily check that $\tilde{S}_W^+(j)/\tilde{S}_W^-(j) = (S_W^+(j)/S_W^-(j))^{1-\rho}$. So, the ratio is (of course) unchanged if $\rho = 0$, and pushed to a pure chance level if $\rho = 1$. We provide below an interpretation of boosting as a greedy optimization procedure that will lead to the value $\rho = 1/2$.

### 10.4.3  Adaboost and greedy gradient descent

We here restrict to the case of binary base classifiers and denote their linear combination by

$$h(x) = \sum_{j=1}^{M} \alpha_j f_j(x).$$

Whether an observation $x$ is correctly classified in the true class $y$ is associated to the sign of the product $yh(x)$, but the value of this product also has an important interpretation, since, when it is positive, it can be thought of as a margin with which $x$ is correctly classified.

Assume that the function $F$ is evaluated, not only on the basis of its classification error, but also based on this margin, using a loss function of the kind

$$\Psi(h) = \sum_{k=1}^{N} \psi(y_k h(x_k)) \tag{10.4}$$

where $\psi$ is decreasing. The boosting algorithm can then be interpreted as an classifier which incrementally improves this objective function.

Let, for $j < M$,

$$h^{(j)} = \sum_{q=1}^{j} \alpha_q f_q.$$

The next combination $h^{(j+1)}$ is equal to $h^{(j)} + \alpha_{j+1} f_{j+1}$, and we now consider the problem of minimizing, with respect to $f_{j+1}$ and $\alpha_{j+1}$, the function $\Psi(h^{(j+1)})$, without modifying the previous classifiers (i.e., performing a greedy optimization). So, we want to minimize, with respect to the base classifier $\tilde{f}$ and to $\alpha \geq 0$, the function

$$U(\alpha, \tilde{f}) = \sum_{k=1}^{N} \psi\left( y_k h^{(j)}(x_k) + \alpha y_k \tilde{f}(x_k) \right)$$

Using the fact that $\tilde{f}$ is a binary classifier, this can be written

$$U(\alpha, \tilde{f}) = \sum_{k=1}^{N} \psi(y_k h^{(j)}(x_k) + \alpha) \mathbf{1}_{y_k = \tilde{f}(x_k)} + \sum_{k=1}^{N} \psi(y_k h^{(j)}(x_k) - \alpha) \mathbf{1}_{y_k \neq \tilde{f}(x_k)} \tag{10.5}$$

$$= \sum_{k=1}^{N} (\psi(y_k h^{(j)}(x_k) - \alpha) - \psi(y_k h^{(j)}(x_k) + \alpha)) \mathbf{1}_{y_k \neq \tilde{f}(x_k)}$$

$$+ \sum_{k=1}^{N} \psi(y_k h^{(j)}(x_k) + \alpha).$$

This shows that $\alpha$ and $\tilde{f}$ have inter-dependent optimality conditions. For a given $\alpha$, the best classifier $\tilde{f}$ must minimize a weighted empirical error with non-negative weights (since $\psi$ is decreasing)

$$w_k = \psi(y_k h^{(j)}(x_k) - \alpha) - \psi(y_k h^{(j)}(x_k) + \alpha).$$

Given $\tilde{f}$, $\alpha$ must minimize the expression in (10.5). One can use an alternative minimization procedure to optimize both $\tilde{f}$ (as a weighted basic classifier) and $\alpha$. However, for the special choice $\psi(t) = e^{-t}$, this optimization turns out to only require one step.

In this case, we have

$$U(\alpha, \tilde{f}) = \sum_{k=1}^{N}(e^{\alpha} - e^{-\alpha})e^{-y_k h^{(j)}(x_k)}\mathbf{1}_{y_k \neq \tilde{f}(x_k)} + e^{-\alpha}\sum_{k=1}^{N}e^{-y_k h^{(j)}(x_k)}$$

$$= e^{-\alpha^{(j)}}(e^{\alpha} - e^{-\alpha})\sum_{k=1}^{N}w_k(j)\mathbf{1}_{y_k \neq \tilde{f}(x_k)} + e^{-\alpha^{(j)}}e^{-\alpha}\sum_{k=1}^{N}w_k(j)$$

with $w_k(j+1) = e^{\alpha^{(j)} - y_k h^{(j)}(x_k)}$ and $\alpha^{(j)} = \alpha_1 + \cdots + \alpha_j$. This shows that $\tilde{f}$ should minimize

$$\sum_{k=1}^{N}w_k(j+1)\mathbf{1}_{y_k \neq \tilde{f}(x_k)}.$$

We note that

$$w_k(j+1) = w_k(j)e^{\alpha_j(1 - y_k f_j(x_k))} = w_k(j)e^{\alpha_j|y_k - f_k(x_k)|},$$

which is identical to the weight updates in algorithm Algorithm 10.4 (this is the reason why the term $\alpha^{(j)}$ was introduced in the computation). The new value of $\alpha$ must minimize (using the notation of Algorithm 10.4)

$$e^{-\alpha}S_W^+(j) + e^{\alpha}S_W^-(j),$$

which yields $\alpha = \frac{1}{2}\log S_W^+(j)/S_W^-(j)$. This is the value $\alpha_{j+1}$ in Algorithm 10.4 with $\rho = 1/2$.

## 10.5  Gradient boosting and regression

### 10.5.1  Notation

The boosting idea, and in particular its interpretation as a greedy gradient procedure, can be extended to non-linear regression problems [75]. Let us denote by $\mathcal{F}_0$

the set of base predictors, therefore functions from $\mathcal{R}_X = \mathbb{R}^d$ to $\mathcal{R}_Y = \mathbb{R}^q$, since we are considering regression problems. The final predictor is a linear combination

$$F(x) = \sum_{j=1}^{M} \alpha_j f_j(x)$$

with $\alpha_1, \ldots, \alpha_M \in \mathbb{R}$ and $f_1, \ldots, f_M \in \mathcal{F}_0$. Note that the the coefficients $\alpha_j$ are redundant when the class $\mathcal{F}_0$ is invariant by multiplication by a scalar. Replacing if needed $\mathcal{F}_0$ by $\{f = \alpha g, \alpha \in \mathbb{R}, g \in \mathcal{F}_0\}$, we will assume that this property holds and therefore remove the coefficients $\alpha_j$ from the problem.

In accordance with the principle of performing greedy searches, we let

$$F^{(j)}(x) = \sum_{q=1}^{j} f_q(x),$$

and consider the problem of minimizing over $f \in \mathcal{F}_0$,

$$U(f) = \sum_{k=1}^{N} r(y_k, F^{(j)}(x_k) + f(x_k)),$$

where $T = (x_1, y_1, \ldots, x_N, y_N)$ is the training data and $r$ is the loss function.

### 10.5.2   Translation-invariant loss

In the case, which is frequent in regression, when $r(y, y')$ only depends on $y - y'$, the problem is equivalent to minimizing

$$U(f) = \sum_{k=1}^{N} r(y_k - F^{(j)}(x_k), f(x_k)),$$

i.e., to let $f_{j+1}$ be the optimal predictor (in $\mathcal{F}_0$ and for the loss $r$) of the residuals $y_k^{(j)} = y_k - F^{(j)}(x_k)$. In this case, this provides a conceptually very simple algorithm.

---

**Algorithm 10.5 (Gradient boosting for regression with translation-invariant loss)**
• Let $T = (x_1, y_1, \ldots, x_N, y_N)$ be a training set and $r$ a loss function such that $r(y, y')$ only depends on $y - y'$.

• Let $\mathcal{F}_0$ be a function class such that $f \in \mathcal{F}_0 \Rightarrow \alpha f \in \mathcal{F}_0$ for all $\alpha \in \mathbb{R}$.

• Select an integer $M > 0$ and let $F^{(0)} = 0$, $y_k^{(0)} = y_k$, $k = 1, \ldots, N$.

• For $j = 1, \ldots, M$:

(1) Find the optimal predictor $f_j \in \mathcal{F}_0$ for the training set $(x_1, y_1^{(j-1)}, \ldots, x_N, y_N^{(j-1)})$.

(2) Let $y_k^{(j)} = y_k^{(j-1)} - f_j(x_k)$

- Return $F = \sum_{k=1}^M f_j$.

---

**Remark 10.6** Obviously, the class $\mathcal{F}_0$ should not be a linear class for the boosting algorithm to have any effect. Indeed, if $f, f' \in \mathcal{F}_0$ implies $f + f' \in \mathcal{F}_0$, no improvement could be made to the predictor after the first step.                                                        ◆

A successful example of this algorithm uses regression trees as base predictors. Recall that the functions output by such trees take the form

$$f(x) = \sum_{A \in \mathcal{C}} w_A \mathbf{1}_{x \in A}$$

where $\mathcal{C}$ is a finite partition of $\mathbb{R}^d$. Each set in the partition is specified by the value taken by a finite number of binary features (denoted by $\gamma$ in our discussion of prediction trees) and the maximal number of such features is the depth of the tree. We assume that the set $\Gamma$ of binary features is shared by all regression trees in $\mathcal{F}_0$, and that the depth of these trees is bounded by a fixed constant. These restrictions prevent $\mathcal{F}_0$ from forming a linear class.[2] Note that the maximal depth of tree learnable from a finite training set is always bounded, since such trees cannot have more nodes than the size of the training set (but one may want to restrict the maximal depth of base predictors to be way less than $N$).

### 10.5.3  General loss functions

We now consider situations in which the loss function is not necessarily a function of the difference between true and predicted output. We are still interested in the problem of minimizing $U(f)$, but we now approximate this problem using the first-order expansion

$$U(f) = \sum_{k=1}^N r(y_k, F^{(j)}(x_k)) + \sum_{k=1}^N \partial_2 r(y_k, F^{(j)}(x_k))^T f(x_k) + o(f),$$

where $\partial_2 r$ denotes the derivative of $r$ with respect to its second variable. This suggests (similarly to gradient descent) to choose $f$ such that $f(x_k) = -\alpha \partial_2 r(y_k, F^{(j)}(x_k))$

---

[2]If $f$ and $g$ are representable as trees, $f + g$ can be represented as a tree whose depth is the sum as those of the original trees, simply by inserting copies of $g$ below each leaf of $f$.

for some $\alpha > 0$ and all $k = 1, \ldots, N$. However, such an $f$ may not exist in the class $\mathcal{F}_0$, and the next best choice is to pick $f = \alpha \tilde{f}$ with $\tilde{f}$ minimizing

$$\sum_{k=1}^{N} |\tilde{f}(x_k) + \partial_2 r(y_k, F^{(j)}(x_k))|^2$$

over all $\tilde{f} \in \mathcal{F}_0$. This is similar to projected gradient descent in optimization, and $\alpha$ such that $f = \alpha \tilde{f}$ should minimize

$$\sum_{k=1}^{N} r(y_k, F^{(j)}(x_k) + \alpha \tilde{f}(x_k)).$$

This provides a generic "gradient boosting" algorithm [75], summarized below.

---

**Algorithm 10.6 (Gradient boosting)**
- Let $T = (x_1, y_1, \ldots, x_N, y_N)$ be a training set and $r$ a differentiable loss function.
- Let $\mathcal{F}_0$ be a function class such that $f \in \mathcal{F}_0 \Rightarrow \alpha f \in \mathcal{F}_0$ for all $\alpha \in \mathbb{R}$.
- Select an integer $M > 0$ and let $F^{(0)} = 0$.
- For $j = 1, \ldots, M$:

  (1) Find $\tilde{f}_j \in \mathcal{F}_0$ minimizing

  $$\sum_{k=1}^{N} |\tilde{f}(x_k) + \partial_2 r(y_k, F^{(j-1)}(x_k))|^2$$

over all $\tilde{f} \in \mathcal{F}_0$.

  (2) Let $f_j = \alpha_j \tilde{f}_j$ where $\alpha_j$ minimizes

  $$\sum_{k=1}^{N} r(y_k, F^{(j-1)}(x_k) + \alpha \tilde{f}_j(x_k)).$$

  (3) Let $F^{(j)} = F^{(j-1)} + f_j$.

- Return $F = F^{(M)}$.

---

**Remark 10.7** Importantly, the fact that $\mathcal{F}_0$ is stable by scalar multiplication implies that the function $\tilde{f}_j$ satisfies

$$\sum_{k=1}^{N} \tilde{f}(x_k)^T \partial_2 r(y_k, F^{(j-1)}(x_k)) \leq 0,$$

$\blacklozenge$

that is, excepted in the unlikely case in which the above sum is zero, it is a direction of descent for the function $U$ (because one could otherwise replace $\tilde{f}_j$ by $-\tilde{f}_j$ and improve the approximation of the gradient).

### 10.5.4   Return to classification

A slight modification of this algorithm may also be applied to classification, provided that the classifier $f$ is obtained by learning the conditional distribution, denoted $g \mapsto p(g|x)$, of the output variable (assumed to take values in a finite set $\mathcal{R}_Y$) given the input (assumed to take values in $\mathcal{R}_X = \mathbb{R}^d$).

Our goal is to estimate an unknown target conditional distribution, $\mu$, therefore taking the form $\mu(g|x)$ for $g \in \mathcal{R}_Y$ and $x \in \mathbb{R}^d$. We assume that a family $\mu_k, k = 1, \ldots, N$ of distributions on the set $\mathcal{R}_Y$ is observed, where each $\mu_k$ is assumed to be an approximation of the unknown $\mu(\cdot|x_k)$ (typically, $\mu_k(g) = \mathbf{1}_{g=y_k}$, i.e., $\mu_k = \delta_{y_k}$). The risk function must take the form $r(\mu, \mu')$ where $\mu, \mu' \in \mathcal{S}(\mathcal{R}_Y)$, the set of probability distributions on $\mathcal{R}_Y$. We will work with

$$r(\mu, \mu') = - \sum_{g \in \mathcal{R}_Y} \mu(g) \log \mu'(g).$$

One can note that

$$r(\mu, \mu') = KL(\mu\|\mu') + r(\mu, \mu),$$

which is therefore minimal when $\mu' = \mu$. Moreover, in the special case $\mu_k = \delta_{y_k}$, the empirical risk is

$$\hat{R}(p) = \sum_{k=1}^{N} r(\mu_k, p(\cdot|x_k)) = - \sum_{k=1}^{N} \log p(y_k|x_k),$$

so that minimizing it is equivalent to maximizing the conditional likelihood that was used for logistic regression.

Before applying the previous algorithm, one must address the issue that probability distributions do not form a vector space, and cannot be added to form new probability distributions. In Friedman [75], Hastie et al. [87], it is suggested to use the representation, which can be associated with any function $F : (g, x) \mapsto F(g|x) \in \mathbb{R}$,

$$p_F(g|x) = \frac{e^{F(g|x)}}{\sum_{h \in \mathcal{R}_Y} e^{F(h|x)}}.$$

Because the representation if not unique ($p_F = p_{F'}$ if $F - F'$ only depends on $x$), we will require in addition that

$$\sum_{h \in \mathcal{R}_Y} F(h|x) = 0$$

for all $x \in \mathbb{R}^d$. The space formed by such functions $F$ is now linear, and we can consider the empirical risk

$$\hat{R}(F) = - \sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} \mu_k(g) \log p_F(g|x_k) = - \sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} \mu_k(g) F(g|x_k) + \sum_{k=1}^{N} \log \left( \sum_{g \in \mathcal{R}_Y} e^{F(g|x_k)} \right).$$

One can evaluate the derivative of this risk with respect to a change on $F(g|x_k)$, and a short computation gives

$$\frac{\partial R}{\partial F(g|x_k)} = -\sum_{k=1}^{N} (\mu_k(g) - p_F(g|x_k)).$$

Now assume that a basic space $\mathcal{F}_0$ of functions $f : (g,x) \mapsto f(g|x)$ is chosen, such that all function in $\mathcal{F}_0$ satisfy

$$\sum_{g \in \mathcal{R}_Y} f(g|x) = 0$$

for all $x \in \mathbb{R}^d$. The gradient boosting algorithm then requires to minimize (in Step (1)):

$$\sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} (\mu_k(g) - p_{F^{(j-1)}}(g|x_k) - \tilde{f}(g|x_k))^2$$

with respect to all functions $\tilde{f} \in \mathcal{F}_0$. Given the optimal $\tilde{f}_j$, the next step requires to minimize, with respect to $\alpha \in \mathbb{R}$:

$$-\alpha \sum_{k=1}^{N} \sum_{g \in \mathcal{R}_Y} \mu_k(g) \tilde{f}_j(g|x_k) + \sum_{k=1}^{N} \log \left( \sum_{g \in \mathcal{R}_Y} e^{F^{(j-1)}(g|x_k) + \alpha \tilde{f}_j(g|x_k)} \right).$$

This is a scalar convex problem that can be solved, e.g., using gradient descent.

### 10.5.5 Gradient tree boosting

We now specialize to the situation in which the set $\mathcal{F}_0$ contains regression trees. In this situation, the general algorithm can be improved by taking advantage of the fact that the predictors returned by such trees are piecewise constant functions, where the regions of constancy are associated with partitions $\mathcal{C}$ of $\mathbb{R}^d$ defined by the leaves of the trees. In particular, $\tilde{f}_j(x)$ in Step (1) takes the form

$$\tilde{f}_j(g|x) = \sum_{A \in \mathcal{C}}^{J} \tilde{f}_{j,A}(g) \mathbf{1}_{x \in A}.$$

The final $f$ at Step (2) should therefore take the form

$$\sum_{A \in \mathcal{C}} \alpha \tilde{f}_{j,A}(g) \mathbf{1}_{x \in A}$$

but not much additional complexity is introduced by freely optimizing the values of $f_j$ on $A$, that is, by looking at $f$ in the form

$$\sum_{A \in \mathcal{C}} f_{j,A}(g) \mathbf{1}_{x \in A}$$

where the values $f_{j,A}(g)$ optimize the empirical risk. This risk becomes

$$-\sum_{k=1}^{N} \sum_{A \in \mathcal{C}} \sum_{g \in \mathcal{R}_Y} \mu_k(g) f_{j,A}(g) \mathbf{1}_{x_k \in A} + \sum_{k=1}^{N} \sum_{A \in \mathcal{C}} \log \left( \sum_{g \in \mathcal{R}_Y} e^{F^{(j-1)}(g|x_k) + f_{j,A}(g)} \right) \mathbf{1}_{x_k \in A}.$$

The values $f_{j,A}(g), g \in \mathcal{R}_Y$ can therefore be optimized separately, minimizing

$$-\sum_{k=1:x_k \in A} \sum_{g \in \mathcal{R}_Y} \mu_k(g) f_{j,A}(g) + \sum_{k:x_k \in A} \log \left( \sum_{g \in \mathcal{R}_Y} e^{F^{(j-1)}(g|x_k) + f_{j,A}(g)} \right) \mathbf{1}_{x_k \in A}.$$

This is still a convex program, which has to be run at every leaf of the optimized tree. If computing time is limited (or for large-scale problems), the determination of $f_{j,A}(g)$ may be restricted to one step of gradient descent starting at $f_{j,A} = 0$. A simple computation indeed shows that the first derivative of the function above with respect to $f_{j,A}(g)$ is

$$a_A(g) = -\sum_{k:x_k \in A} (\mu_k(g) - p_F(g|x_k)).$$

The derivative of this expression with respect to $f_{j,A}(g)$ (for the same $g$) is

$$b_A(g) = \sum_{k:x_k \in A} p_F(g|x_k)(1 - p_F(g|x_k)).$$

The off-diagonal terms in the second derivative are, for $g \neq h$,

$$-\sum_{k:x_k \in A} p_F(g|x_k) p_F(h|x_k).$$

In Friedman et al. [74], it is suggested to use an approximate Newton step, where the off-diagonal terms in the second derivative are neglected. This corresponds to minimizing

$$\sum_{g \in \mathcal{R}_Y} a_A(g) f_{j,A}(g) + \frac{1}{2} \sum_{g \in \mathcal{R}_Y} b_A(g) f_{j,A}(g)^2.$$

The solution is (introducing a Lagrange multiplier for the constraint $\sum_g f_{j,A}(g) = 0$)

$$f_{j,A}(g) = -\frac{a_A(g) - \lambda}{b_A(g)}$$

with

$$\lambda = \frac{\sum_{g \in \mathcal{R}_Y} a_A(g)/b_A(g)}{\sum_{g \in \mathcal{R}_Y} 1/b_A(g)}.$$

A small value $\epsilon$ can be added to $b_A$ to avoid divisions by zero. We refer the reader to Friedman et al. [74], Friedman [75], Hastie et al. [87] for several variations on this basic idea. Note that an approximate but highly efficient implementation of boosted trees, called XGBoost, has been developed in Chen and Guestrin [53].

# Chapter 11

# Iterated Compositions of Functions and Neural Nets

## 11.1   First definitions

We now discuss a class of methods in which the predictor $f$ is built using iterated compositions, with a main application to neural nets. We will structure these models using directed acyclic graphs (DAG). These graphs are composed with a set of vertexes (or nodes) $\mathcal{V} = \{0, \ldots, m+1\}$ and a collection $\mathcal{L}$ of directed edges $i \to j$ between some vertexes. If an edge exists between $i$ and $j$, one says that $i$ is a parent of $j$ and $j$ a child of $i$ and will use the notation $pa(i)$ (resp. $ch(i)$) denote the set of parents (resp. children) of $i$. The graphs we consider must satisfy the following conditions:

   (i)  No index is a descendant of itself, i.e., that the graph is acyclic.

   (ii)  The only index without parent is $i = 0$ and the only one without children in $i = m+1$.

   To each node $i$ in the graph, one associates a dimension $d_i$ and a variable $z_i \in \mathbb{R}^{d_i}$. The root node variable, $z_0 = x$, is the input and $z_{m+1}$ is the output. One also associates to each node $i \neq 0$ a function $\psi_i$ defined on the product space $\bigotimes_{j \in pa(i)} \mathbb{R}^{d_j}$ and taking values in $\mathbb{R}^{d_i}$. The input-output relation is then defined by the family of equations:

$$z_i = \psi_i(z_{pa(i)})$$

where $z_{pa(i)} = (z_j, j \in pa(i))$. Since there is only one root and one terminal node, these iterations implement a relationship $y = z_{m+1} = f(x)$, with $z_0 = x$. We will refer to the $z_1, \ldots, z_m$ as the *latent variables* of the network.

   Each function $\psi_i$ is furthermore parametrized by an $s_i$-dimensional vector $w_i \in \mathbb{R}^{s_i}$, so that we will write

$$z_i = \psi_i(z_{pa(i)}; w_i).$$

We let $\mathcal{W}$ denote the vector containing all parameters $w_1, \ldots, w_{m+1}$, which therefore has dimension $s = s_1 + \cdots + s_{m+1}$. The network function $f$ is then parametrized by $\mathcal{W}$ and we will write $y = f(x; \mathcal{W})$.

## 11.2   Neural nets

### 11.2.1   Transitions

Most neural networks iterate functions taking the form

$$\psi_i(z; w) = \rho(bz + \beta_0), z \in \mathbb{R}^{d_j}$$

where $b$ is a $d_i \times (\sum_{j \in pa(i)} d_j)$ matrix and $\beta_0 \in \mathbb{R}^{d_i}$ (so that $w = (b, \beta_0)$ is $s_i = d_i(1 + \sum_{j \in pa(i)} d_j)$-dimensional); $\rho$ is defined on and takes values in $\mathbb{R}$, and we make the abuse of notation, for any $d$ and $u \in \mathbb{R}^d$

$$\rho(u) = \begin{pmatrix} \rho(u^{(1)}) \\ \vdots \\ \rho(u^{(d)}) \end{pmatrix}.$$

The most popular choice for $\rho$ is the positive part, or ReLU (for rectified linear unit), given by $\rho(t) = \max(t, 0)$. Other common choices are $\rho(t) = 1/(1 + e^{-t})$ (sigmoid function), or $\rho(t) = \tanh(z)$.

Residual neural networks (or ResNets [89]) are discussed in section 11.6. They iterate transitions between inputs and outputs of same dimension, taking

$$z_{i+1} = z_i + \psi(z_i; w). \tag{11.1}$$

### 11.2.2   Output

The last node of the graph provides the prediction, $y$. Its expression depends on the type of predictor that is learned

• For regression, $y$ can be chosen as an affine function of is its parents: $z_{m+1} = bz_{pa(m+1)} + a_0$.

• For classification, one can also use a linear model $z_{m+1} = bz_{pa(m+1)} + a_0$ where $z_{m+1}$ is $q$-dimensional and let the classification be $\mathrm{argmax}(z_{m+1}^{(i)}, i = 1, \ldots, q)$. Alternatively, one uses "softmax" transformation, with

$$z_{m+1}^{(i)} = \frac{e^{\zeta_{m+1}^{(i)}}}{\sum_{j=1}^q e^{\zeta_{m+1}^{(j)}}}$$

with $\zeta_{m+1} = bz_{pa(m+1)} + a_0$.

### 11.2.3 Image data

Neural networks have achieved top performance when working with organized structures such as images. A typical problem in this setting is to categorize the content of the image, i.e., return a categorical variable naming its principal element(s). Other applications include facial recognition or identification. In this case, the transition function can take advantage of the 2D structure, with some special terminology.

Instead of speaking of the total dimension, say, $d$, of the considered variables, writing $z = (z^{(1)}, \ldots, z^{(d)})$, images are better represented with three indices $z(u, v, \lambda)$ where $u = 1, \ldots, U$ and $U$ is the width of the image, $v = 1, \ldots, V$ and $V$ is the height of the image, $\lambda = 1, \ldots, \Lambda$ and $\Lambda$ is the depth of the image. (With this notation $d = UV\Lambda$.) Typical images have length and width of one or two hundred pixels, and depth $\Lambda = 3$ for the three color channels. This three-dimensional structure is conserved also for latent variables, with different dimensions. Deep neural networks often combine compression in width and height with expansion in depth while transitioning from input to output.

The linear transformation $b$ mapping one layer with dimensions $U_i, V_i, \Lambda_i$ to another with dimensions $U_{i+1}, V_{i+1}, \Lambda_{i+1}$ is then preferably seen as a collection of numbers: $b(u', v', \lambda', u, v, \Lambda)$ so that the transition from $z_i$ to $z_{i+1}$ is given by

$$z_{i+1}(u', v', \lambda') = \rho \left( \beta_0(u', v', \lambda') + \sum_{u=1}^{U_i} \sum_{v=1}^{V_i} \sum_{\lambda=1}^{\Lambda_i} b(u', v', \lambda', u, v, \lambda) z_i(u, v, \lambda) \right).$$

For images, it is often preferable to use convolutional transitions, providing convolutional neural networks ([116, 115], or CNNs. If $U_i = U_{i+1}$ and $V_i = V_{i+1}$, such a transition requires that $b(u', v', \lambda', u, v, \lambda)$ only depends on $\lambda$, $\lambda'$ and on the differences $u' - u$ and $v - v'$. In general, one also requires that $b(u', v', \lambda', u, v, \lambda)$ is non-zero only if $|u' - u|$ and $|v' - v|$ are both less than a constant, typically a small number. Also, there is generally little computation across depths: each output at depth $\lambda'$ only uses values from a single input depth. These restrictions obviously reduce dramatically the number of free parameters involved in the transition.

After one or a few convolutions, the dimension is often reduced by a "pooling" operation, dividing the image into small non-overlapping windows and replacing each such window by a single value, either the max (max-pooling) or the average.

## 11.3 Geometry

In addition to the transitions between latent variables and resulting changes of dimension, the structure of the DAG defining the network is an important element in

Figure 11.1: Linear net with increasing layer depths and decreasing layer width.



Figure 11.2: A sketch of the U-net architecture designed for image segmentation [168].

the design of a neural net. The simplest choice is a purely linear structure (as shown in Figure 11.1), as was, for example, used for image categorization in [110].

More complex architectures have been introduced in recent years. Their design is in a large part heuristic and based on an analysis of the kind of computation that should be done in the network to perform a particular task. For example, an architecture used for image segmentation in summarized in fig. 11.2.

An important feature of neural nets is their modularity, since "simple" architectures can be combined (e.g., by placing the output of a network as input of another one) and form a more complex network that still follows the basic structure defined above. One example of such a building block is the "attention module," which take as input three vectors $Q, K, V$ (for query, key, and value) and return

$$\text{softmax}(QK^T)V.$$

These modules are fundamental elements of "transformer networks" [198], that are used, among other tasks, for automatic translation.

## 11.4 Objective function

### 11.4.1 Definitions

We now return to the general form of the problem, with variables $z_0, \ldots, z_{m+1}$ satisfying

$$z_i = \psi_i(z_{pa(i)}; w_i)$$

Let $T = (x_1, y_1, \ldots, x_N, y_N)$ denote the training data.

For regression problems, the objective function minimized by the algorithm is typically the empirical risk, the simplest choice being the mean square error, which gives

$$F(\mathcal{W}) = \frac{1}{N} \sum_{k=1}^{N} |y_k - z_{k,m+1}(\mathcal{W})|^2.$$

with $z_{k;m+1}(\mathcal{W}) = f(x_k; \mathcal{W})$.

For classification, with the dimension of the output variable equal to the number of classes and the decision based on the largest coordinate, one can take (letting $z_{k,m+1}(i; \mathcal{W})$ denote the $i$th coordinate of $z_{k,m+1}(\mathcal{W})$):

$$F(\mathcal{W}) = \frac{1}{N} \sum_{k=1}^{N} \left( -z_{k,m+1}(y_k; \mathcal{W}) + \log \left( \sum_{i=1}^{q} \exp(z_{k,m+1}(i; \mathcal{W})) \right) \right).$$

This objective function is similar to that minimized in logistic regression.

### 11.4.2 Differential

**General computation.** The computation of the differential of $F$ with respect to $\mathcal{W}$ may look daunting, but it has actually a simple structure captured by the back-propagation algorithm. Even if programming this algorithm can often be avoided by using an automatic differentiation software, it is important to understand how it works, and why the implementation of gradient-descent algorithms remains feasible.

Consider the general situation of minimizing a function $G(\mathcal{W}, z)$ over $\mathcal{W} \in \mathbb{R}^s$ and $z \in \mathbb{R}^r$, subject to a constraint $\gamma(\mathcal{W}, z) = 0$ where $\gamma$ is defined on $\mathbb{R}^s \times \mathbb{R}^r$ and takes values in $\mathbb{R}^r$ (here, it is important that the number of constraints is equal to the dimension of $z$). We will denote below by $\partial_{\mathcal{W}}$ and $\partial_z$ the derivatives of these functions with respect to the multi-dimensional variables $\mathcal{W}$ and $z$. We make the assumptions that $\partial_z \gamma$, which is an $r \times r$ matrix, is invertible, and that the constraints can be solved to express $z$ as a function of $\mathcal{W}$, that we will denote $Z(\mathcal{W})$.

This allows us to define the function $F(\mathcal{W}) = G(\mathcal{W}, \boldsymbol{Z}(\mathcal{W}))$ and we want to compute the gradient of $F$. (Clearly, the function $F$ in the previous section satisfies these assumptions). Taking $h \in \mathbb{R}^s$, we have

$$dF(\mathcal{W})h = \partial_{\mathcal{W}} G h + \partial_{\boldsymbol{z}} G \, d\boldsymbol{Z} h.$$

Moreover, since $\gamma(\mathcal{W}, \boldsymbol{Z}) = 0$ by definition of $\boldsymbol{Z}$, we have

$$\partial_{\mathcal{W}} \gamma h + \partial_{\boldsymbol{z}} \gamma \, d\boldsymbol{Z} h = 0,$$

so that

$$dF(\mathcal{W})h = \partial_{\mathcal{W}} G h - \partial_{\boldsymbol{z}} G \, \partial_{\boldsymbol{z}} \gamma^{-1} \, \partial_{\mathcal{W}} \gamma h.$$

Let $\boldsymbol{p} \in \mathbb{R}^r$ be the solution of the linear system

$$\partial_{\boldsymbol{z}} \gamma^T \boldsymbol{p} = \partial_{\boldsymbol{z}} G^T.$$

Then,

$$dF(\mathcal{W})h = (\partial_{\mathcal{W}} G - \boldsymbol{p}^T \partial_{\mathcal{W}} \gamma)h$$

or

$$\nabla F = \partial_{\mathcal{W}} G^T - \partial_{\mathcal{W}} \gamma^T \boldsymbol{p}.$$

Note that, introducing the "Hamiltonian"

$$\boldsymbol{H}(\boldsymbol{p}, \boldsymbol{z}, \mathcal{W}) = \boldsymbol{p}^T \gamma(\mathcal{W}, \boldsymbol{z}) - G(\mathcal{W}, \boldsymbol{z}),$$

one can summarize the previous computation with the system

$$\begin{cases} \partial_{\boldsymbol{p}} \boldsymbol{H} = 0 \\ \partial_{\boldsymbol{z}} \boldsymbol{H} = 0 \\ \nabla F = -\partial_{\mathcal{W}} \boldsymbol{H}^T. \end{cases}$$

**Application: back-propagation.**   In our case, we are minimizing a function of the form

$$G(\mathcal{W}, \boldsymbol{z}_1, \ldots, \boldsymbol{z}_N) = \frac{1}{N} \sum_{k=1}^{N} r(y_k, z_{k,m+1})$$

subject to constraints $z_{k,i+1} = \psi_i(z_{k,pa(i)}; w_i)$, $i = 0, \ldots, m$, $z_{k,0} = x_k$. We focus on one of the terms in the sum, therefore fixing $k$, that we will temporarily drop from the notation.

So, we evaluate the gradient of $G(\mathcal{W}, z) = r(y, z_{m+1})$ with $z_{i+1} = \psi_i(z_{pa(i)}; w_i)$, $i = 0, \ldots, m$, $z_0 = x$. With the notation of the previous paragraph, we take $\gamma = (\gamma_1, \ldots, \gamma_{m+1})$ with

$$\gamma_i(\mathcal{W}, z) = \psi_i(z_{pa(i)}; w_i) - z_i$$

These constraints uniquely define $z$ as a function of $\mathcal{W}$, which was one of our assumptions. For the derivative, we have, for $u = (u_1, \ldots, u_{m+1}) \in \mathbb{R}^r$ (with $r = d_1 + \cdots + d_{m+1}$, $u_i \in \mathbb{R}^{d_i}$), and for $i = 1, \ldots, m+1$

$$\partial_z \gamma_i u = \sum_{j \in pa(i)} \partial_{z_j} \psi_i(z_{pa(i)}; w_i) u_j - u_i$$

Taking $p = (p_1, \ldots, p_{m+1}) \in \mathbb{R}^r$, we get

$$p^T \partial_z \gamma u = \sum_{i=1}^{m+1} \sum_{j \in pa(i)} p_i^T \partial_{z_j} \psi_i(z_{pa(i)}; w_i) u_j - \sum_{i=1}^{m+1} p_i^T u_i$$

$$= \sum_{j=1}^{m+1} \sum_{i \in ch(j)} p_i^T \partial_{z_j} \psi_i(z_{pa(i)}; w_i) u_j - \sum_{j=1}^{m+1} p_j^T u_j$$

This allows us to identify $\partial_z \gamma^T p$ as the vector $g = (g_1, \ldots, g_{m+1})$ with

$$g_j = \sum_{i \in ch(j)} \partial_{z_j} \psi_i(z_{pa(i)}; w_i)^T p_i - p_j.$$

For $j = m + 1$ (which has no children), we get $g_{m+1} = -p_{m+1}$, so that the equation $\partial_z \gamma^T p = g$ can be solved recursively by taking $p_{m+1} = -g_{m+1}$ and propagating backward, with

$$p_j = -g_j + \sum_{i \in ch(j)} \partial_{z_j} \psi_i(z_{pa(i)}; w_i)^T p_i$$

for $j = m, \ldots, 1$.

To compute the gradient of $G$, the propagation has to be applied to $g = \partial_z G$. Since $G$ only depends on $z_{m+1}$, we have $g_{m+1} = \partial_{z_{m+1}} r(y, z_{m+1})$ and $g_j = 0$ for $j = 1, \ldots, m$. Moreover, $G$ does not depend on $\mathcal{W}$, so that $\partial_{\mathcal{W}} G = 0$. We have

$$\partial_{\mathcal{W}} \gamma_i = \partial_{w_i} \psi_i(z_{pa(i)}, w_i)$$

yielding $\partial_{\mathcal{W}} \gamma^T p = (\zeta_1, \ldots, \zeta_m)$ with

$$\zeta_j = \partial_{w_j} \psi_j(z_{pa(j)}, w_j)^T p_j.$$

We can now formulate an algorithm that computes the gradient of $F$ with respect to $\mathcal{W}$, reintroducing training data indexes in the notation.

**Algorithm 11.1 (Back-propagation)**
Let $(x_1, y_1, \ldots, x_N, y_N)$ be the training set and $R_k(z) = r(y_k, z)$ so that

$$F(\mathcal{W}) = \frac{1}{N} \sum_{k=1}^{N} R_k(z_{k,m+1}(\mathcal{W}))$$

with $z_{k,m+1}(\mathcal{W}) = f(x_k, \mathcal{W})$. Let $\mathcal{W}$ be a family of weights. The following steps compute $\nabla F(\mathcal{W})$.

1. For all $k = 1, \ldots, N$ and all $i = 1, \ldots, m+1$, compute $z_{k,i}(\mathcal{W})$ (forward computation through the network).

2. Initialize variables $p_{k,m+1} = -\nabla R_k(z_{k,m+1}(\mathcal{W}))$, $k = 1, \ldots, N$.

3. For all $k = 1, \ldots, N$ and all $j = 1, \ldots, m$, compute $p_{k,j}$ using iterations

$$p_{k,j} = \sum_{i \in ch(j)} \partial_{z_j} \psi_i(z_{k,pa(i)}, w_i)^T p_{k,i}.$$

4. Let

$$\nabla F(\mathcal{W}) = -\frac{1}{N} \sum_{k=1}^{N} \sum_{i=1}^{m+1} D_i^T \partial_{w_i} \psi_i(z_{k,pa(i)}, w_i)^T p_{k,i},$$

where $D_i$ is the $s_i \times s$ matrix such that $D_i h = h_i$.

---

### 11.4.3   Complementary computations

The back-propagation algorithm requires the computation of the gradient of the costs $R_k$ and of the derivatives of the functions $\psi_i$, and this can generally be done in closed form, with relatively simple expressions.

- If $R_k(z) = |y_k - z|^2$ (which is the typical choice for regression models) then $\nabla R_k(z) = 2(z - y_k)$.

- In classification, with $R_k(z) = -z(y_k) + \log\left(\sum_{i=1}^{q} \exp(z^{(i)})\right)$, one has

$$\nabla R_k(z) = -u_{y_k} + \frac{\exp(z)}{\sum_{i=1}^{q} \exp(z^{(i)})}$$

where $u_{y_k} \in \mathbb{R}^d$ is the vector with 1 at position $y_k$ and zero elsewhere, and $\exp(z)$ is the vector with coordinates $\exp(z^{(i)})$, $i = 1, \ldots, d$.

- For dense transition functions in the form $\psi(z; w) = \rho(bz + \beta_0)$ with $w = (\beta_0, b)$, then $\partial_z \psi(z, w) = \mathrm{diag}(\rho'(\beta_0 + bz))b$ so that

$$\partial_z \psi(z, w)^T p = b^T \mathrm{diag}(\rho'(\beta_0 + bz))p$$

- Similarly

$$\partial_w \psi(z, w)^T p = \left[\mathrm{diag}(\rho'(\beta_0 + bz))p, \mathrm{diag}(\rho'(\beta_0 + bz))pz^T\right].$$

Note that neural network packages implement these functions (and more) automatically.

## 11.5 Stochastic Gradient Descent

### 11.5.1 Mini-batches

Fix $\ell \ll N$. Consider the set of $B_\ell$ of binary sequences $\xi = (\xi^1, \ldots, \xi^N)$ such that $\xi^k \in \{0, 1\}$ and $\sum_{k=1}^N \xi^k = \ell$. Define

$$H(\mathcal{W}, \xi) = \nabla_{\mathcal{W}} \left( \frac{1}{\ell} \sum_{k=1}^N \xi^k r(y_k, f(x_k, \mathcal{W})) \right) = \frac{1}{\ell} \sum_{k=1}^N \xi^k \nabla_{\mathcal{W}} r(y_k, f(x_k, \mathcal{W}))$$

where $\xi$ follows the uniform distribution on $B_\ell$. Consider the stochastic approximation algorithm:

$$\mathcal{W}_{n+1} = \mathcal{W}_n - \gamma_{n+1} H(\mathcal{W}_n, \xi_{n+1}). \tag{11.2}$$

Because $E(\xi^k) = \ell/N$, we have $E(H(\mathcal{W}, \xi)) = \nabla_{\mathcal{W}} \mathcal{E}_T(f(\cdot, \mathcal{W}))$ and (11.2) provides a stochastic gradient descent algorithm to which the discussion in section 3.3 applies. Such an approach is often referred to as "mini-batch" selection in the deep-learning literature, since it correspond to sampling $\ell$ examples from the training set without replacement and only computing the gradient of the empirical loss computed from these examples.

### 11.5.2 Dropout

Introduced for deep learning in Srivastava et al. [181], "dropout" is a learning paradigm that brings additional robustness (and, maybe, reduces overfitting risks) to massively parametrized predictors.

Assume that a random perturbation mechanism of the model parameters has been designed. We will represent it using a random variable $\eta$ (interpreted as noise) and a transformation $\mathcal{W}' = \varphi(\mathcal{W}, \eta)$ describing how $\eta$ affects a given weight configuration $\mathcal{W}$ to form a perturbed one $\mathcal{W}'$. In order to shorten notation, we will

write $\varphi(\mathcal{W}, \eta) = \eta \cdot \mathcal{W}$, borrowing the notation for a group action from group theory. As a typical example, $\eta$ can be chosen as a vector of Bernoulli random variables (therefore taking values in $\{0,1\}$), with same dimension as $\mathcal{W}$ and one can simply let $\eta \cdot \mathcal{W} = \eta \odot \mathcal{W}$ be the pointwise multiplication of the two vectors. This corresponds to replacing some of the parameters by zero ("dropping them out") while keeping the others unchanged. One generally preserves the parameters of the final layer ($g_m$), so that the corresponding $\eta$'s are equal to one, and let the other ones be independent, with some probability $p$ of being one, say, $p = 1/2$.

Returning to the general case, in which $\eta$ is simply assumed to be a random variable with known probability distribution, the dropout method replaces the objective function $F(\mathcal{W}) = \mathcal{E}_T(f(\cdot, \mathcal{W}))$ by its expectation over perturbed predictors $G(\mathcal{W}) = E(\mathcal{E}_T(f(\cdot, \eta \cdot \mathcal{W})))$ where the expectation is taken with respect to the random variable $\eta$. While this expectation cannot be computed explicitly, its minimization can be performed using stochastic gradient descent, with

$$\mathcal{W}_{n+1} = \mathcal{W}_n - \gamma_{n+1} L(\mathcal{W}_n, \eta_{n+1}),$$

where $\eta_1, \eta_2, \dots$ is a sequence of independent realizations of $\eta$ and

$$L(\mathcal{W}, \eta) = \nabla_\mathcal{W} \left( \mathcal{E}_T(f(\cdot, \eta \cdot \mathcal{W})) \right).$$

Then, averaging in $\eta$

$$\bar{L}(\mathcal{W}) = E(\nabla_\mathcal{W} F(\eta \odot \mathcal{W})) = \nabla G(\mathcal{W}).$$

In the special case where $\eta \cdot \mathcal{W}$ is just pointwise multiplication, then

$$L(\mathcal{W}, \eta) = \eta \odot \nabla F(\eta \odot \mathcal{W}).$$

So this quantity can be evaluated by using back-propagation to compute $\nabla F(\eta \cdot \mathcal{W})$ and multiplying the result by $\eta$ pointwise. Obviously, random weight perturbation can be combined with mini-batch selection in a hybrid stochastic gradient descent algorithm, the specification of which being left to the reader. We also note that stochastic gradient descent in neural networks is often implemented using the ADAM algorithm (section 3.3.3).

## 11.6   Continuous time limit and dynamical systems

### 11.6.1   Neural ODEs

Equation (11.1) expresses the difference of the input and output of a neural transition as a non-linear function $f(z; w)$ of the input. This strongly suggests passing to

continuous time and replacing the difference by a derivative, i.e., replacing the neural network by a high-dimensional parametrized dynamical system. The continuous model then takes the form [52]

$$\partial_t z(t) = \psi(z(t); w(t)) \tag{11.3}$$

where $t$ varies in a a fixed interval, say, $[0, T]$. The whole process is parametrized by $\mathcal{W} = (w(t), t \in [0, T])$. We need to assume existence and uniqueness of solutions of (11.3), which usually restricts the domain of admissibility of parameters $\mathcal{W}$.

Typical neural transition functions are Lipschitz functions whose constant depend on the weight magnitude, i.e., are such that

$$|\psi(z, w) - \psi(z', w)| \le C(w)|z - z'| \tag{11.4}$$

where $C$ is a continuous function of $W$. For example, for $\psi(z, w) = \rho(bz + \beta_0)$, $w = (b, \beta_0)$, one can take $C(w) = C_\rho |b|_{\text{op}}$. The Caratheodory theorem [17] implies that solutions are well-defined as soon as

$$\int_0^T C(w(t))dt < \infty. \tag{11.5}$$

This is a relatively mild requirement, on which we will return later. Assuming this, we can consider $z(T)$ as a function of the initial value, $z(0) = x$ and of the parameters, writing $z(T) = f(x, \mathcal{W})$.

Given a training set, we consider the problem of minimizing

$$F(\mathcal{W}) = \frac{1}{N} \sum_{k=1}^N r(y_k, f(x_k, \mathcal{W})). \tag{11.6}$$

The discussion in section section 11.4.2 applies—formally, at least—to this continuous case, and we can consider the equivalent problem of minimizing

$$G(\mathcal{W}, z_1, \ldots, z_N) = \frac{1}{N} \sum_{k=1}^N r(y_k, z_k(T))$$

with $\partial_t z_k(t) = \psi(z_k(t); w(t))$, $z_k(0) = x_k$. Once again, we consider each $k$ separately, which boils down to considering $N = 1$ and we drop the index $k$ from the notation, letting $F(\mathcal{W}) = r(y, f(x, \mathcal{W}))$ $G(\mathcal{W}, \bar{z}) = r(y, z(T))$.

We define $\gamma(\mathcal{W}, z)$ to return the *function*

$$t \mapsto \gamma(\mathcal{W}, z)(t) = \psi(z(t); w(t)) - \partial_t z(t).$$

Let $\boldsymbol{p} : [0,T] \to \mathbb{R}^d$.  We want to determine the expression of $\boldsymbol{u} = \partial_z \gamma^T \boldsymbol{p}$, which satisfies

$$\int_0^T \boldsymbol{u}(t)^T \delta \boldsymbol{z}(t) dt = \int_0^T \boldsymbol{p}(t)^T (\partial_z \psi(\boldsymbol{z}(t), w(t)) \delta \boldsymbol{z}(t) - \partial_t \delta \boldsymbol{z}(t)) dt$$

After an integration by parts, the r.h.s. becomes

$$-\boldsymbol{p}(T)^T \delta \boldsymbol{z}(T) + \int_0^T \partial_t \boldsymbol{p}(t)^T \delta \boldsymbol{z}(t) dt + \int_0^T \boldsymbol{p}(t)^T \partial_z \psi(\boldsymbol{z}(t), w(t)) \delta \boldsymbol{z}(t)) dt$$

which gives

$$\boldsymbol{u}(t) = -\boldsymbol{p}(T) \delta_T + \partial_t \boldsymbol{p}(t) + \partial_z \psi(\boldsymbol{z}(t), w(t))^T \boldsymbol{p}(t).$$

The equation $\partial_z \gamma^T \boldsymbol{p} = \partial_z \boldsymbol{G}^T$ therefore gives

$$-\boldsymbol{p}(T) \delta_T + \partial_t \boldsymbol{p}(t) + \partial_z \psi(\boldsymbol{z}(t), w(t))^T \boldsymbol{p}(t) = \partial_2 r(y, \boldsymbol{z}(T)) \delta_T,$$

so that $\boldsymbol{p}$ satisfies $\boldsymbol{p}(T) = -\partial_2 r(y, \boldsymbol{z}(T))$ and

$$\partial_t \boldsymbol{p}(t) = -\partial_z \psi(\boldsymbol{z}(t), w(t))^T \boldsymbol{p}(t). \tag{11.7}$$

We have $\partial_{\mathcal{W}} \boldsymbol{G} = 0$ and $\boldsymbol{v} = \partial_{\mathcal{W}} \gamma^T \boldsymbol{p}$ satisfies

$$\int_0^T \boldsymbol{v}(t)^T \delta w(t) dt = \int_0^T \boldsymbol{p}(t)^T \partial_w \psi(\boldsymbol{z}(t), w(t)) \delta w(t) dt$$

so that

$$\nabla F(\mathcal{W}) = (t \mapsto -\partial_w \psi(\boldsymbol{z}(t), w(t))^T \boldsymbol{p}(t)).$$

This informal derivation (more work is needed to justify the existence of various differentials in appropriate function spaces) provides the continuous-time version of the back-propagation algorithm, which is also known as the *adjoint method* in the optimal control literature [91, 123].  In that context, $z$ represents the state of the control system, $w$ is the control and $p$ is called the costate, or covector.  We summarize the gradient computation algorithm, reintroducing $N$ training samples.

---

**Algorithm 11.2 (Adjoint method for neural ODE)**
Let $(x_1, y_1, \ldots, x_N, y_N)$ be the training set and $R_k(z) = r(y_k, z)$ so that

$$F(\mathcal{W}) = \frac{1}{N} \sum_{k=1}^N R_k(\boldsymbol{z}_k(T, \mathcal{W}))$$

with $\partial_t \boldsymbol{z}_k = \psi(\boldsymbol{z}_k, \mathcal{W})$, $\boldsymbol{z}_k(0) = x_k$.  Let $\mathcal{W}$ be a family of weights.  The following steps compute $\nabla F(\mathcal{W})$.

1. For all $k = 1, \ldots, N$ and all $t \in [0, T]$, compute $z_k(t, \mathcal{W})$ (forward computation through the dynamical system).

2. Initialize variables $p_k(T) = -\nabla R_k(z_k(T, \mathcal{W}))/N$, $k = 1, \ldots, N$.

3. For all $k = 1, \ldots, N$ and all $j = 1, \ldots, m$, compute $p_k(t)$ by solving (backwards in time)
$$\partial_t p_k(t) = -\partial_z \psi(z_k(t), w(t))^T p_k(t).$$

4. Let, for $t \in [0, T]$,
$$\nabla F(\mathcal{W})(t) = -\sum_{k=1}^{N} \partial_w \psi(z_k(t), w(t))^T p_k(t).$$

---

Of course, in numerical applications, the forward and backward dynamical systems need to be discretized, in time, resulting in a finite number of computation steps. This can be done explicitly (for example using basic Euler schemes), or using ODE solvers [52] available in every numerical software.

### 11.6.2 Adding a running cost

Optimal control problems are usually formulated with a "running cost" that penalizes the magnitude of the control, which in our case is provided by the function $\mathcal{W} : t \mapsto w(t)$. Penalties on network weights are rarely imposed with discrete neural networks, but, as discussed above, in the continuous setting, some assumptions on the function $\mathcal{W}$, such as (11.5), are needed to ensure that the problem is well defined.

It is therefore natural to modify the objective function in (11.6) by adding a penalty term ensuring the finiteness of the integral in (11.5), taking, for example, for some $\lambda > 0$,
$$F(\mathcal{W}) = \lambda \int_0^T C(w(t))^2 dt + \sum_{k=1}^{N} r(y_k, f(x_k, \mathcal{W})). \tag{11.8}$$

The finiteness of the integral of the squared $C(w)^2$ implies, by Cauchy-Schwartz, the integrability of $C(w)$ itself, and usually leads to simpler computations.

If $C(w)$ is known explicitly and is differentiable, the previous discussion and the back-propagation algorithm can be adapted with minor modifications for the minimization of (11.8). The only difference appears in Step 4 of Algorithm 11.2, with
$$\nabla F(\mathcal{W})(t) = 2\lambda \nabla C(w(t)) - \frac{1}{N} \sum_{k=1}^{N} \partial_w \psi(z_k(t), w(t))^T p_k(t).$$

Computationally, one should still ensure that $C$ and its gradient are not too costly to compute. If $\psi(z, w) = \rho(bz + \beta_0)$, $w = (b, \beta_0)$, the choice $C(w) = C_\rho |b|_{\text{op}}$ is valid, but not computationally friendly. The simpler choice $C(w) = C_\rho |b|_2$ is also valid, but cruder as an upper-bound of the Lipschitz constant. It leads however to straightforward computations.

The addition of a running cost to the objective is important to ensure that any potential solution of the problem leads to a solvable ODE. It does not guarantee that an optimal solution exists, which is a trickier issue in the continuous setting than in the discrete setting. This is an important theoretical issue, since it is needed, for example, to ensure that various numerical discretization schemes lead to consistent approximations of a limit continuous problem. The existence of minimizers is not known in general for ODE networks. It does hold, however, in the following non-parametric (i.e., weight-free) context that we now describe.

The function $\psi$ in the r.h.s. of (11.3), is, for any fixed $w$, a function that maps $z \in \mathbb{R}^d$ to a vector $\psi(z, w) \in \mathbb{R}^d$. Such functions are called *vector fields* on $\mathbb{R}^d$, and the collection $\psi(\cdot, w), w \in \mathbb{R}^s$ is a parametrized family of vector fields.

The non-parametric approach replaces this family of functions by a general vector field, $v$ so that the time-indexed parametrized family of vector fields $(t \mapsto \psi(\cdot, w(t)))$ becomes an unconstrained family $(t \mapsto f(t, \cdot))$. Following the general non-parametric framework in statistics, one needs to define a suitable function space for the vector fields, and use a penalty in the objective function.

We will assume that, at each time, $f(t, \cdot)$ belongs to a reproducing kernel Hilbert space (RKHS), as introduced in chapter 6. However, because we are considering a space of vector fields rather than scalar-valued functions, we need work with matrix-valued kernels [5], for which we give a definition that generalizes definition 6.1 (which corresponds to $q = 1$ below).

**Definition 11.1** *A function $K : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathcal{M}_q(\mathbb{R})$ satisfying*

*[K1-vec]  K is symmetric, namely $K(x, y) = K(y, x)^T$ for all $x$ and $y$ in $\mathbb{R}^d$.*

*[K2-vec]  For any $n > 0$, for any choice of vectors $\lambda_1, \ldots, \lambda_n \in \mathbb{R}^q$ and any $x_1, \ldots, x_n \in \mathbb{R}^d$, one has*

$$\sum_{i,j=1}^{n} \lambda_i^T K(x_i, x_j) \lambda_j \geq 0. \tag{11.9}$$

*is called a positive (matrix-valued) kernel.*

*One says that the kernel is positive definite if the sum in (6.1) cannot vanish, unless (i) $\lambda_1 = \cdots = \lambda_n = 0$ or (ii) $x_i = x_j$ for some $i \neq j$.*

If $\kappa$ is a "scalar kernel" (satisfying definition 6.1), then $K(x,y) = \kappa(x,y)\mathrm{Id}_{\mathbb{R}^q}$ is a matrix-valued kernel.

A reproducing kernel Hilbert space of vector-valued functions is a Hilbert space $H$ of functions from $\mathbb{R}^d$ to $\mathbb{R}^q$ such that there exists a reproducing kernel $K : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathcal{M}_q(\mathbb{R})$ with the following properties

[RKHS1]  For all $x \in \mathbb{R}^d$ and $\lambda \in \mathbb{R}^q$, $K(\cdot,x)\lambda$ belongs to $H$,

[RKHS2]  For all $h \in H$, $x \in \mathbb{R}^d$ and $\lambda \in \mathbb{R}^q$,

$$\langle h, K(\cdot,x)\lambda \rangle_H = \lambda^T h(x).$$

Proposition 6.5 remains valid in the for vector-valued RKHS, with the following modifications: $\lambda_1,\ldots,\lambda_N$ and $\alpha_1,\ldots,\alpha_N$ are $q$-dimensional vectors and the matrix $\mathcal{K}(x_1,\ldots,x_N)$ is now an $Nq \times Nq$ block matrix, with $q \times q$ blocks given by $K(x_k,x_l)$, $k,l = 1,\ldots,N$.

Returning to the specification of the nonparametric control problem, we will assume that a vector-valued RKHS, $H$, has been chosen, with $q = d$ in definition 11.1. We further assume that elements of $H$ are Lipschitz continuous, with

$$|v(z) - v(\tilde{z})| \leq C\|v\|_H |z - \tilde{z}| \tag{11.10}$$

for some constant $C$ and all $v \in H$. We note that, for every $\lambda \in \mathbb{R}^d$,

$$
\begin{aligned}
|\lambda^T(v(z) - v(\tilde{z}))|^2 &= |\langle v, K(\cdot,z)\lambda - K(\cdot,\tilde{z})\lambda \rangle_H|^2 \\
&\leq \|v\|_H^2 \|K(\cdot,z)\lambda - K(\cdot,\tilde{z})\lambda\|_H^2 \\
&= \|v\|_H^2 (\lambda^T K(z,z)\lambda - 2\lambda^T K(z,\tilde{z})\lambda + \lambda^T K(\tilde{z},\tilde{z})) \\
&\leq |\lambda|^2 \|v\|_H^2 |K(z,z) - 2K(z,\tilde{z}) + K(\tilde{z},\tilde{z})|.
\end{aligned}
$$

This shows that (11.10) can be derived from regularity properties of the kernel, namely, that

$$|K(z,z) - 2K(z,\tilde{z}) + K(\tilde{z},\tilde{z})| \leq C|z - \tilde{z}|^2$$

for some constant $C$ and all $z,\tilde{z} \in \mathbb{R}^d$. This property is satisfied by most of the kernels that are used in practice.

Let $\eta : t \mapsto \eta(t)$ be a function from $[0,1]$ to $H$. This means that, for each $t$, $\eta(t)$ is a vector field $x \mapsto \eta(t)(x)$ on $\mathbb{R}^d$, and we will write indifferently $\eta(t)$ and $\eta(t,\cdot)$, with a preference for $\eta(t,x)$ rather than $\eta(t)(x)$. We consider the objective function

$$\bar{F}(f) = \lambda \int_0^T \|\eta(t)\|_H^2 dt + \frac{1}{N}\sum_{k=1}^N r(y_k, z_k(1)), \tag{11.11}$$

with $\partial_t z_k(t) = \eta(t, z_k(t))$, $z_k(0) = x_k$. To compare with (11.8), the finite-dimensional $w \in \mathbb{R}^s$ is now replaced with an infinite-dimensional parameter, $\eta$, and the transition $\psi(z, w)$ becomes $\eta(z)$.

Using the vector version of proposition 6.5 (or the kernel trick used several times in chapters 7 and 8), one sees that there is no loss of generality in replacing $\eta(t)$ by its projection onto the vector space

$$V(t) = \left\{ \sum_{l=1}^{N} K(\cdot, z_l(t)) w_l : w_1, \ldots, w_N \in \mathbb{R}^d \right\}.$$

Noting that, if $\eta(t)$ takes the form

$$\eta(t) = \sum_{l=1}^{N} K(\cdot, z_l(t)) w_l(t),$$

then

$$\|\eta(t)\|_H^2 = \sum_{k,l=1}^{N} w_k(t)^T K(z_k(t), z_l(t)) w_l(t).$$

This allows us to replace the infinite-dimensional parameter $\eta$ by a family $\mathcal{W} = (w(t), t \in [0, T]$ with $w(t) = (w_k(t), k = 1, \ldots, N)$. The minimization of $\bar{F}$ in (11.11) can be replaced by that of

$$F(\mathcal{W}) = \lambda \int_0^T \sum_{k,l=1}^{N} w_k(t)^T K(z_k(t), z_l(t)) w_l(t) dt + \frac{1}{N} \sum_{k=1}^{N} r(y_k, z_k(1)), \qquad (11.12)$$

with

$$\partial_t z_k(t) = \sum_{l=1}^{N} K(z_k(t), z_l(t)) w_l(t).$$

This optimal control problem has a similar form to that considered in (11.8), where the running cost $C(w)^2$ is replaced by a cost that depends on the control (still denoted $w$) and the state $z$. The discussion in section section 11.6.1 can be applied with some modifications. Let $\mathcal{K}(z)$ be the $dN \times dN$ matrix formed with $d \times d$ blocks $K(z_k(t), z_l(t))$ and $w(t)$ the $dN$-dimensional vector formed by stacking $w_1, \ldots, w_N$. Let

$$G(\mathcal{W}, z) = \lambda \int_0^T w(t)^T \mathcal{K}(z(t)) w(t) dt + \frac{1}{N} \sum_{k=1}^{N} r(y_k, z_k(1))$$

and

$$\gamma(\mathcal{W}, z)(t) = \mathcal{K}(z(t)) w(t) - \partial_t z(t).$$

The backward ODE in step 3. of Algorithm 11.2 now becomes

$$\partial_t \boldsymbol{p}_k(t) = -\partial_{z_k}(\boldsymbol{w}(t)^T \mathcal{K}(\boldsymbol{z}(t))\boldsymbol{p}(t)) + \lambda \partial_{z_k}(\boldsymbol{w}(t)^T \mathcal{K}(\boldsymbol{z}(t))\boldsymbol{w}(t))$$

for $k = 1,\ldots,N$. Step 4. becomes (for $t \in [0,T]$),

$$\nabla F(\mathcal{W})(t) = \mathcal{K}(\boldsymbol{z}(t))(2\lambda - \boldsymbol{p}(t)).$$

The resulting algorithm was introduced in [209]. It has the interesting property (shared with neural ODE models with smooth controlled transitions) to determine an implicit diffeomorphic transformation of the space, i.e., the function $x \mapsto f(x;\mathcal{W},\boldsymbol{z}) = \tilde{z}(T)$ which returns the solution at time $T$ of the ODE

$$\partial_t \tilde{z}(t) = \sum_{l=1}^{N} K(\tilde{z}(t), z_l(t))w_l(t)$$

(or $\partial \tilde{z}(t) = \psi(\tilde{z}(t);w(t))$ for neural ODEs) is smooth, invertible, with a smooth inverse.

# Chapter 12

# Monte-Carlo Sampling

The goal of this section is to describe how, from a basic random number generator that provides samples from a uniform distribution on $[0, 1]$, one can generate samples that follow, or approximately follow, complex probability distributions on finite or general spaces. This, combined with the law of large numbers, permits to approximate probabilities or expectations by empirical averages over a large collection of generated samples.

We assume that as many as needed independent samples of the uniform distribution are available, which is only an approximation of the truth. In practice, computer programs are only able to generate pseudo-random numbers, which are highly chaotic recursive sequences, but still deterministic. Also, these numbers are generated as integers, which only provide, after normalization, a distribution on a finite discretization of the unit interval. We will neglect these facts, however, and work as if the output of the function *random* (or any similar name) in a computer program is a true realization of the uniform distribution.

## 12.1 General sampling procedures

**Real-valued variables.** We will use the following notation for the left limit of a function $F$ at a given point $z$

$$F(z - 0) = \lim_{y \to z, y < z} F(y)$$

assuming, of course that this limit exists (which is always true, for example when $F$ is non-decreasing). Recall that $F$ is left continuous if and only if $F = F(\cdot - 0)$. Moreover, it is easy to see that $F(\cdot - 0)$ is left-continuous[1]. Note also that, if $F$ is non-decreasing,

---

[1]For every $z$ and every $\epsilon > 0$, there exists $z' < z$ such that for all $z'' \in [z', z)$, $|F(z'') - F(z - 0)| < \epsilon$. Moreover, taking any $y \in (z', z)$, there exists $y' < y$ such that for all $y'' \in [y', y)$, $|F(y'') - F(y - 0)| < \epsilon$.

one always has $F(z) \leq F(y-0)$ whenever $z < y$. The following proposition provides a basic mechanism for Monte-Carlo sampling.

**Proposition 12.1** *Let $Z$ be a real-valued random variable with c.d.f. $F_Z$. For $u \in [0,1]$, define*

$$F_X^-(u) = \max\{z : F_Z(z-0) \leq u\}.$$

*Let $U$ be uniformly distributed over $[0,1]$. Then $F_Z^-(U)$ has the same distribution as $Z$.*

PROOF  Let $A_z = \{u \in [0,1] : F_Z^-(u) \leq z\}$. Assume first that $u < F_Z(z)$. Then $F_Z(y-0) \leq u$ implies that $y \leq z$, since $y > z$ would imply that $F_Z(z) \leq F_Z(y-0)$. This shows that $\sup\{z' : F_Z(z'-0) \leq u\} \leq z$, i.e., $u \in A_z$.

Now, take $u > F_Z(z)$. Because c.d.f.'s are right continuous, there exists $y > z$ such that $u > F_Z(y)$, which implies that $F_Z^-(u) \geq y$ and $u \notin A_z$.

We have therefore shown that $[0, F_Z(z)) \subset A_z \subset [0, F_Z(z)]$. If $U$ is uniformly distributed on $[0,1]$, then $P(U < F_Z(z)) = P(U \leq F_Z(z)) = F_Z(z)$, showing that

$$\mathbb{P}(F_Z^-(U) \leq z) = \mathbb{P}(U \in A_z) = F_Z(z). \qquad \blacksquare$$

This proposition shows that one can generate random samples of a real-valued random variable $Z$ as soon as one can compute $F_Z^-$ and generate uniformly distributed variables. Note that, if $F_Z$ is strictly increasing, then $F_Z^- = F_Z^{-1}$, the usual function inverse.

The proposition also shows how to sample from random variables taking values in finite sets. Indeed, if $Z$ takes values in $\widetilde{\Omega}_Z = \{z_1, \ldots, z_n\}$ with $p_i = \mathbb{P}(Z = z_i)$, sampling from $Z$ is equivalent to sampling from the integer valued random variable $\widetilde{Z}$ with $\mathbb{P}(\widetilde{Z} = i) = p_i$. For this variable, $F_{\widetilde{Z}}^-(u)$ is the largest $i$ such that $p_1 + \cdots + p_{i-1} \leq u$ (this sum being zero if $i = 1$), which provides the standard sampling scheme for discrete probability distributions.

## 12.2  Rejection sampling

While the previous approach can be generalized to multivariate distributions, it quickly becomes unfeasible when the dimension gets large, excepting simple cases in which the variables are independent, or, say, Gaussian. Rejection sampling is a simple algorithm that allows, in some cases, for the generation of samples from a complicated distribution based on repeated sampling of a simpler one.

---

Without loss of generality, we can assume that $y' \geq z'$, yielding $|F(z-0) - F(y-0)| \leq 2\epsilon$, showing the left continuity of $F(\cdot - 0)$.

Let us assume that we want to sample from a variable $Z$ taking values $\mathcal{R}_Z$, and that there exists a measure $\mu$ on $\mathcal{R}_Z$ with respect to which the distribution of $Z$ is absolutely continuous, i.e., so that this distribution has a density $f_Z$ with respect to $\mu$. For example, $\mathcal{R}_Z = \mathbb{R}^d$, and $f_Z$ is the p.d.f. of $Z$ with respect to Lebesgue's measire. Assume that $g$ is another density functions (with respect to $\mu$) from which it is "easy" to sample. Consider the following algorithm, which includes a function $a : z \mapsto a(z) \in [0,1]$ that will be specified later.

---

**Algorithm 12.1 (Rejection sampling with acceptance function $a$ and base p.d.f. $g$)**
(1) Sample a realization $z$ of a random variable with p.d.f. $g$.

(2) Generate $b \in \{0,1\}$ with $\mathbb{P}(b = 1) = a(z)$.

(3) If $b = 1$, return $Z = z$ and exit.

(4) Otherwise, return to step 1.

---

The probability of exiting at step 3 is $\rho = \int_{\mathbb{R}^d} g(z)a(z)\mu(dz)$. So, the algorithm simulates a random variable with p.d.f.

$$\tilde{f}(z) = g(z)a(z)(1 + (1 - \rho) + (1 - \rho)^2 + \cdots) = \frac{g(z)a(z)}{\rho}.$$

As a consequence, in order to simulate $f_Z$, one must choose $a$ so that $f_Z(z)$ is proportional to $g(z)a(z)$, which, (assuming that $g(z) > 0$ whenever $f_Z(z) > 0$), requires that $a(z)$ is proportional to $f_Z(z)/g(z)$. Since $a(z)$ must take values in $[0,1]$, but should otherwise be chosen as large as possible to ensure that fewer iterations are needed, one should take

$$a(z) = \frac{f_Z(z)}{cg(z)}$$

where $c = \max\{f_Z(z)/g(z) : z \in \mathbb{R}^d\}$, which must therefore be finite. This fully specifies a rejection sampling algorithm for $f_Z$. Note that $g$ is free to choose (with the restriction that $f_Z(z)/g(z)$ must be bounded), and should be selected so that sampling from it is easy, and the coefficient $c$ above is not too large.

## 12.3 Markov chain sampling

When dealing with high-dimensional distributions, the constant $c$ in the previous procedure is typically extremely large, and the rejection-sampling algorithm becomes unfeasible, because it keeps rejecting samples for very long times. In such cases, one can use alternative simulation methods that iteratively updates the variable $Z$ by making small changes at each step, resulting in a procedure that asymptotically converges to a sample of the target distribution. Such sampling schemes

are usually described as Markov chains, leading to the name Markov-chain Monte Carlo (or MCMC) sampling.

Assume that we want to sample from a random variable that takes values in some (measurable) set $\mathcal{B} = \mathcal{R}_X$.[2] A Markov chain is the probabilistic analogous of a recursive sequence $X_{n+1} = \Phi(X_n)$, which is fully defined by the function $\Phi : \mathcal{B} \to \mathcal{B}$ and the initial value $X_0 \in \mathcal{B}$.

### 12.3.1   Definitions

For Markov chains, $X_0$ is a random variable, which therefore does not have a fixed value, but follows a probability distribution that we will generally denote $\mu_0$: $P^0(x) = P(X_0 = x)$. The computation of $X_{n+1}$ given $X_n$ is not deterministic either, but given the conditional probabilities

$$P^{n,n+1}(x, A) = \mathbb{P}(X_{n+1} \in A \mid X_n = x).$$

where $A \subset \mathcal{B}$ is measurable. The left-hand side of this equation, $P^{n,n+1}$ is called a transition probability, according to the following definition.

**Definition 12.2** *Let $F_1$ and $F_2$ be two sets equipped with $\sigma$-algebras $\mathcal{A}_1$ and $\mathcal{A}_2$. A transition probability from $F_1$ to $F_2$ is a function $p : F_1 \times \mathcal{A}_2 \to [0, 1]$ such that, for all $x \in F_1$, the function $A \mapsto p(x, A)$ is a probability on $F_2$ and for all $A \in \mathcal{A}_2$, the function $x \mapsto p(x, A)$, $x \in F_1$, is measurable.*

*When $F_2$ is discrete, the probabilities are fully specified by their values on singleton sets, and we will write $p(x, y)$ for $p(x, \{y\})$.*

When $P^{n,n+1}(x, \cdot)$ does not depend on $n$, the Markov chain is said to be homogeneous. To simplify notation, we will restrict to homogeneous chains (and therefore only write $P(x, A)$), although some of the chains used in MCMC sampling may be inhomogeneous. This is not a very strong loss of generality, however, because inhomogeneous Markov chains can be considered as homogeneous by extending the space $\Omega$ on which they are defined to $\Omega \times \mathbb{N}$, and defining the transition probability

$$\tilde{p}\big((x, n), A \times \{r\}\big) = \mathbf{1}_{r=n+1} p^{n,n+1}(x, A).$$

An important special case is when $\mathcal{B}$ is countable, in which case one only needs to specify transition probabilities for singletons $A = \{y\}$, and we will write

$$p(x, y) = P(x, \{y\}) = \mathbb{P}(X_{n+1} = y \mid X_n = x)$$

---

[2]We will assume in this chapter that $\mathcal{B}$ is a complete metric space with a dense countable subset, with the associated Borel $\sigma$-algebra.

for the p.m.f. associated with $P(x,\cdot)$.

Another simple situation is when $\mathcal{B} = \mathbb{R}^d$ and each $P(x,\cdot)$ has a p.d.f. that we will also denote as $p(x,\cdot)$. In this latter case, assuming that $P^0$ also have a p.d.f. that we will denote by $\mu_0$, the joint p.d.f. of $(X_0,\ldots,X_n)$ on $(\mathbb{R}^d)^{n+1}$ is given by

$$f(x_0,x_1,\ldots,x_n) = \mu_0(x_0)p(x_0,x_1)\cdots p(x_{n-1},x_n). \tag{12.1}$$

The same expression holds for the joint p.m.f. in the discrete case.

In the general case (invoking measure theory), the joint distribution is also determined by the transition probabilities, and we leave the derivation of the expression to the reader. An important point is that, in both special cases considered above, and under some very mild assumptions in the general case , these transition probabilities also uniquely define the joint distribution of the infinite process $(X_0,X_1,\ldots)$ on $\mathcal{B}^\infty$, which gives theoretical support to the consideration of asymptotic properties of Markov chains. In this discussion, we are interested in conditions ensuring that the chain asymptotically samples from a target probability distribution $Q$, i.e., whether $\mathbb{P}(X_n \in A)$ converges to $Q(A)$ (one says that $X_n$ converges in distribution to $Q$). In practice, $Q$ is given or modeled, and the goal is to determine the transition probabilities. Note that the marginal distribution of $X_n$ is computed by integrating (or summing) (12.1) with respect to $x_0,\ldots,x_{n-1}$, which is generally computationally challenging.

Given a transition probability $P$ on $\mathcal{B}$, we will use the notation, for a function $f : \mathcal{B} \to \mathbb{R}$:

$$Pf(x) = \int_{\mathcal{B}} f(y)P(x,dy).$$

If $Q$ is a probability distribution on $\mathcal{B}$, it will also be convenient to write

$$Qf(x) = \int_{\mathcal{B}} f(y)Q(dy).$$

### 12.3.2 Convergence

We will denote $\mathbb{P}_x(\cdot)$ the conditional distribution $\mathbb{P}(\cdot \mid X_0 = x)$ and $P^n(x,A) = \mathbb{P}_x(X_n \in A)$, which is a probability distribution on $\mathcal{B}$. The goal of Markov Chain Monte Carlo sampling is to design the transition probabilities such that $P_x^n(A)$ converges to $Q(A)$ when $n$ tends to infinity. One furthermore wants to complete this convergence with a law of large numbers, ensuring that

$$\frac{1}{n}\sum_{k=1}^{n} f(X_k) \to \int_{\mathcal{B}} f(x)Q(dx)$$

when $n \to \infty$, where $X_n$ is the generated Markov chain and $f$ is $Q$-integrable.

Introduce the total variation distance between two probability measures on a given probability space,

$$D_{\text{var}}(\mu_1, \mu_2) = \sup_A (\mu_1(A) - \mu_2(A)). \tag{12.2}$$

where the supremum is taken over all measurable sets $A$. We will say that the Markov chain with transition $P$ asymptotically samples from $Q$ if

$$\lim_{n \to \infty} D_{\text{var}}(P^n(x, \cdot), Q) = 0 \tag{12.3}$$

for $Q$-almost all $x \in \mathcal{B}$. The chain must satisfy specific conditions for this to be guaranteed.

We now discuss some properties of the total variation distance that will be useful later. First, we note that the supremum in the r.h.s. of (12.2) is achieved. Indeed, there exists a set $A_0$ such that, for all measurable sets $A$, $\mu_1(A \cap A_0) \geq \mu_2(A \cap A_0)$ and $\mu_1(A \cap A_0^c) \leq \mu_2(A \cap A_0^c)$. If $\mathcal{B}$ is a finite set, it suffices to let $A_0 = \{x \in \mathcal{B} : \mu_1(x) \geq \mu_2(x)\}$; if both $\mu_1$ and $\mu_2$ have p.d.f.'s $\varphi_1$ and $\varphi_2$ with respect to Lebesgue's measure (with $\mathcal{B} = \mathbb{R}^d$), then one can take $A_0 = \{x \in \mathcal{B} : \varphi_1(x) \geq \varphi_2(x)\}$. In the general case, one can take $\mu = \mu_1 + \mu_2$ so that $\mu_1, \mu_2 \ll \mu$, and letting $\varphi_i = d\mu_i/d\mu$, also take $A_0 = \{x \in \mathcal{B} : \varphi_1(x) \geq \varphi_2(x)\}$. (This is also a special case of the Hahn-Jordan decomposition of signed measures [66]).

Now, it is clear that, for any $A$

$$
\begin{aligned}
\mu_1(A) - \mu_2(A) &= \mu_1(A \cap A_0) - \mu_2(A \cap A_0) + \mu_1(A \cap A_0^c) - \mu_2(A \cap A_0^c) \\
&\leq \mu_1(A \cap A_0) - \mu_2(A \cap A_0) \\
&\leq \mu_1(A \cap A_0) - \mu_2(A \cap A_0) + \mu_1(A^c \cap A_0) - \mu_2(A^c \cap A_0) \\
&= \mu_1(A_0) - \mu_2(A_0)
\end{aligned}
$$

showing that

$$D_{\text{var}}(\mu_1, \mu_2) = \mu_1(A_0) - \mu_2(A_0).$$

The following proposition lists additional properties.

**Proposition 12.3** *(i) If $\mu_1, \mu_2$ have a densities $\varphi_1, \varphi_2$ with respect to some positive measure $\mu$ (such as $\mu_1 + \mu_2$), then*

$$D_{\text{var}}(\mu_1, \mu_2) = \frac{1}{2} \int_{\mathcal{B}} |\varphi_1(x) - \varphi_2(x)| \mu(dx).$$

*In particular, if $\mathcal{B}$ is finite*

$$D_{\text{var}}(\mu_1, \mu_2) = \frac{1}{2} \sum_{x \in \mathcal{B}} |\mu_1(x) - \mu_2(x)|.$$

*(ii) For general $\mathcal{B}$,*

$$D_{\text{var}}(\mu_1, \mu_2) = \sup_f \left( \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) \right). \tag{12.4}$$

*where the supremum is taken over all measurable functions $f$ taking values in $[0,1]$.*

*(iii) If $f : \mathcal{B} \to \mathbb{R}$ is bounded, define the maximal oscillation of $f$ by*

$$osc(f) = \sup\{f(x) - f(y) : x, y \in \mathcal{B}\}.$$

*Then*

$$D_{\text{var}}(\mu_1, \mu_2) = \sup \left\{ \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) : osc(f) \le 1 \right\}$$

*(iv) Conversely, for any bounded measurable $f : \mathcal{B} \to \mathbb{R}$,*

$$osc(f) = \sup \left\{ \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) : D_{\text{var}}(\mu_1, \mu_2) \le 1 \right\}$$

PROOF If one takes $A_0 = \{x \in \mathcal{B} : \varphi_1(x) \ge \varphi_2(x)\}$, then

$$D_{\text{var}}(\mu_1, \mu_2) = \int_{A_0} (\varphi_1(x) - \varphi_2(x))\mu(dx) = \int_{A_0} |\varphi_1(x) - \varphi_2(x)|\mu(dx).$$

But, because both $\mu_1$ and $\mu_2$ are probability measures

$$\int_{\mathcal{B}} (\varphi_1(x) - \varphi_2(x))\mu(dx) = 0$$

so that

$$\int_{A_0^c} (\varphi_1(x) - \varphi_2(x))\mu(dx) = - \int_{A_0} (\varphi_1(x) - \varphi_2(x))\mu(dx).$$

However, the l.h.s. is also equal to

$$- \int_{A_0^c} |\varphi_1(x) - \varphi_2(x)|\mu(dx)$$

so that

$$\int_{\mathcal{B}} |\varphi_1(x) - \varphi_2(x)|\mu(dx) = 2 \int_{A_0} (\varphi_1(x) - \varphi_2(x)) = 2D_{\text{var}}(\mu_1, \mu_2),$$

which proves (i).

To prove (ii), first notice that, for all $A$,

$$\mu_1(A) - \mu_2(A) = \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx)$$

for $f = \mathbf{1}_A$, so that

$$D_{\mathrm{var}}(\mu_1, \mu_2) \leq \sup_f \left( \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) \right).$$

Conversely, using $A_0$ as above, and taking $f$ with values in $[0,1]$

$$\int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) = \int_{A_0} f(x)(\mu_1 - \mu_2)(dx) + \int_{A_0^c} f(x)(\mu_1 - \mu_2)(dx)$$

$$\leq \int_{A_0} f(x)(\mu_1 - \mu_2)(dx)$$

$$\leq \int_{A_0} (\mu_1 - \mu_2)(dx) = D_{\mathrm{var}}(\mu_1, \mu_2)$$

This shows (ii). For (iii), one can note that, if $f$ takes values in $[0,1]$, then $osc(f) \leq 1$ so that

$$D_{\mathrm{var}}(\mu_1, \mu_2) \leq \sup \left\{ \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) : osc(f) \leq 1 \right\}$$

Conversely, take $f$ such that $osc(f) \leq 1$, $\epsilon > 0$ and $y$ such that $f(y) \geq inf_x f(x) + \epsilon$. Let $f_\epsilon(x) = (f(x) - f(y) + \epsilon)/(1 + \epsilon)$, which takes values in $[0,1]$. Then

$$D_{\mathrm{var}}(\mu_1, \mu_2) \geq \int_{\mathcal{B}} f_\epsilon(x)\mu_1(dx) - \int_{\mathcal{B}} f_\epsilon(x)\mu_2(dx)$$

$$= \frac{1}{1 + \epsilon} \left( \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) \right)$$

and since this is true for all $\epsilon > 0$, we get

$$\int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) \leq D_{\mathrm{var}}(\mu_1, \mu_2)$$

which completes the proof of (iii).

Using (iii), we find, for any $\mu_1, \mu_2$ and any bounded $f$

$$\int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) \leq D_{\mathrm{var}}(\mu_1, \mu_2)osc(f)$$

which shows that

$$\sup \left\{ \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) : D_{\mathrm{var}}(\mu_1, \mu_2) \leq 1 \right\} \leq osc(f).$$

However, taking $\mu_1 = \delta_x$ and $\mu_2 = \delta_y$, so that $D_{\text{var}}(\mu_1, \mu_2) = 0$ is $x = y$ and 1 otherwise, we get

$$f(x) - f(y) = \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx)$$

$$\leq \sup\left\{ \int_{\mathcal{B}} f(x)\mu_1(dx) - \int_{\mathcal{B}} f(x)\mu_2(dx) : D_{\text{var}}(\mu_1, \mu_2) \leq 1 \right\}$$

which yields (iv) after taking the supremum with respect to $x$ and $y$. ∎

### 12.3.3 Invariance and reversibility

If a Markov chain converges to $Q$, then $Q$ must be an "invariant distribution," in the sense that, if $X_n \sim Q$ for some $n$, then so does $X_{n+1}$ and as a consequences all $X_m$ for $n \geq m$. This can be seen by writing

$$P^{n+1}(x, A) = \mathbb{P}_x(X_{n+1} \in A) = \mathbb{E}_x(P(X_n, A)) = E_{P^n(x, \cdot)}(P(\cdot, A))$$

If $P^n(x, \cdot)$ (and therefore also $P^{n+1}(x, \cdot)$) converges to $Q$, then passing to the limit above yields

$$Q(A) = E_Q(P(\cdot, A))$$

and this states that, if $X_n \sim Q$, then so does $X_{n+1}$. If $Q$ has a p.d.f. (resp. p.m.f.), say, $q$, this gives

$$q(y) = \int_{\mathbb{R}^d} p(x, y)q(x)dx, (\text{resp. } q(y) = \sum_{x \in \mathcal{B}} p(x, y)q(x)).$$

So, if one designs a Markov chain with a target asymptotic distribution $Q$, the first thing to ensure is that $Q$ is invariant.

While invariance leads to an integral equation for $q$, a stronger condition, called reversibility is easier to assess.

Assume that $Q$ is invariant by $P$. Make the assumption that $P(x, \cdot)$ has a density $p_*$ with respect to $Q$ (this is, essentially, no loss of generality, see argument below), so that

$$P(x, A) = \int_A p_*(x, y)Q(dy).$$

Taking $A = \mathcal{B}$ above, we have

$$\int_{\mathcal{B}} p_*(x, y)Q(dy) = P(x, \mathcal{B}) = 1$$

but we also have, because $Q$ is invariant, that

$$\int_{\mathcal{B}} p_*(x, y)Q(dx) = Q(\mathcal{B}) = 1.$$

One says that the density is "doubly stochastic" with respect to $Q$.

Conversely, if a transition probability $P$ has a doubly stochastic density $p_*$ with respect to some probability $Q$ on $\mathcal{B}$, then $Q$ is invariant by $P$, since

$$\int_\mathcal{B} P(x,A)Q(dx) = \int_\mathcal{B} \int_A p_*(x,y)Q(dy)Q(dx)$$

$$= \int_A \int_\mathcal{B} p_*(x,y)Q(dx)Q(dy) = \int_A Q(dy) = Q(A).$$

The property of being doubly stochastic can be reinterpreted in terms of time reversal for Markov chains. Let $Q_0$ be an initial distribution for a Markov chain with transition $P$ (not necessarily invariant) so that, for any $n \geq 0$, the distribution of $X_n$ is $Q_n = Q_0 P^n$. Fixing any $m > 0$, we are interested in the reversed process $\tilde{X}_k = X_{m-k}$. We first notice that the conditional distribution of $X_n$ given its future $X_{n+1}, \ldots, X_m$ (with $n < m$) only depends on $X_{n+1}$, so that the reversed process is also Markov. Indeed, for any positive function $f : \mathcal{B} \to \mathbb{R}$, $g : \mathcal{B}^{m-n} \to \mathbb{R}$, one has, using the fundamental properties of conditional expectations and the fact that $(X_n)$ is a Markov chain,

$$\mathbb{E}(f(X_n)g(X_{n+1},\ldots,X_m)) = \mathbb{E}\left(\mathbb{E}(f(X_n)g(X_{n+1},\ldots,X_m) \mid X_n, X_{n+1})\right)$$
$$= \mathbb{E}\left(f(X_n)\mathbb{E}(g(X_{n+1},\ldots,X_m) \mid X_n, X_{n+1})\right)$$
$$= \mathbb{E}\left(f(X_n)\mathbb{E}(g(X_{n+1},\ldots,X_m) \mid X_{n+1})\right)$$
$$= \mathbb{E}\left(E(f(X_n) \mid X_{n+1})\mathbb{E}(g(X_{n+1},\ldots,X_m) \mid X_{n+1})\right)$$
$$= \mathbb{E}\left(E(f(X_n) \mid X_{n+1})g(X_{n+1},\ldots,X_m)\right).$$

This shows that

$$\mathbb{E}(f(X_n) \mid X_{n+1},\ldots,X_m) = \mathbb{E}(f(X_n) \mid X_{n+1}),$$

which is what we wanted. To identify the conditional distribution of $X_n$ given $X_{n+1}$, we note that for any $x \in \mathcal{B}$, the transition probability $P(x,\cdot)$ is absolutely continuous with respect to $Q_{n+1}$ since

$$Q_{n+1}(A) = \int_\mathcal{B} P(x,A)Q_n(dx)$$

and the r.h.s. is zero only if $P(x,A) = 0$ $Q_n$-almost everywhere [3]. This shows that there exists a function $r_{n+1} : \mathcal{B} \times \mathcal{B} \to \mathbb{R}$ such that, for all $A$

$$P(x,A) = \int_A r_{n+1}(x,y)Q_{n+1}(dy).$$

---

[3]The "almost everywhere" statement a priori depends on $A$, but can be made independent of it under the mild assumption (that we will always make) that $\mathcal{B}$ has a countable basis of open sets.

Given this point, one can write

$$\mathbb{E}(f(X_n)g(X_{n+1})) = \int_{\mathcal{B}^2} f(x_n)g(x_{n+1})P(x_n, dx_{n+1})Q_n(dx_n)$$

$$= \int_{\mathcal{B}^2} f(x_n)g(x_{n+1})r_{n+1}(x_n, x_{n+1})Q_{n+1}(dx_{n+1})Q_n(dx_n)$$

$$= \int_{\mathcal{B}} \left( \int_{\mathcal{B}} f(x_n)r_{n+1}(x_n, x_{n+1})Q_n(dx_n) \right) g(x_{n+1})Q_{n+1}(dx_{n+1})$$

which shows that the conditional distribution of $X_n$ given $X_{n+1} = x_{n+1}$ has density $x_n \mapsto r_{n+1}(x_n, x_{n+1})$ relatively to $Q_n$. Note that, for discrete probabilities, one has

$$r_{n+1}(x, y) = \frac{P(x, y)}{Q_{n+1}(y)}$$

and

$$\mathbb{P}(X_n = x \mid X_{n+1} = y) = \frac{Q_n(x)P(x, y)}{Q_{n+1}(y)}. \tag{12.5}$$

The formula is identical if both $Q_0$ and $P(x, \cdot)$ have p.d.f.'s with respect to a fixed reference measure $\mu$ on $\mathcal{B}$ (for example, Lebesgue's measure when $\mathcal{B} = \mathbb{R}^d$), denoting these p.d.f's by $q_0$ and $p(x, \cdot)$. Then, the p.d.f. of the distribution of $X_n$ given $X_{n+1} = y$ is

$$\tilde{p}_n(y, x) = \frac{q_n(x)p(x, y)}{q_{n+1}(y)} \tag{12.6}$$

where $q_n$ is the p.d.f. of $Q_n$. Note that the transition probabilities of the reversed Markov chain depend on $n$, i.e., the reversed chain is non-homogeneous in general.

However, if one assumes that $Q_0 = Q$ is invariant by $P$, then $Q_n = Q$ for all $n$ and therefore $r_n(x, y) = p_*(x, y)$, using the previous notation. In this case, the reversed chain has transitions independent of $n$ and its transition probability has density

$$\tilde{p}_*(x, y) = p_*(y, x)$$

with respect to $Q$. In the discrete case, letting $p(x, y) = P(X_{n+1} = y \mid X_n = x)$, we have $p_*(x, y) = p(x, y)/Q(y)$, so that the reversed transition (call it $\tilde{p}$) is such that

$$\frac{\tilde{p}(x, y)}{Q(y)} = \frac{p(y, x)}{Q(x)},$$

i.e.,

$$Q(y)p(y, x) = Q(x)\tilde{p}(x, y). \tag{12.7}$$

One retrieves easily the fact that if $p$ is such that there exists $Q$ and $\tilde{p}$ such that (12.7) is satisfied, then (summing the equation over $y$) $Q$ is an invariant probability for $p$.

Let $Q$ be a probability on $\mathcal{B}$. One says that the Markov chain (or the transition probability $p$) is $Q$-reversible if and only if $p(x,\cdot)$ has a density $p_*(x,\cdot)$ with respect to $Q$ such that $p_*(x,y) = p_*(y,x)$ for all $x,y \in \mathcal{B}$. Since such a density is necessarily doubly stochastic, $Q$ is then invariant by $p$. Reversibility is equivalent to the property that, whenever $X_n \sim Q$, the joint distribution of $(X_n, X_{n+1})$ coincides with that of $(X_{n+1}, X_n)$. Alternatively, $Q$-reversibility requires that for all $A, B \subset \mathcal{B}$,

$$\int_A P(z,B)dQ(z) = \int_B P(z,A)dQ(z). \tag{12.8}$$

In the discrete case, (12.8) is equivalent to the "detailed balance" condition:

$$Q(y)p(y,x) = Q(x)p(x,y). \tag{12.9}$$

While $Q$ can be an invariant distribution for a Markov chain without that chain being $Q$-reversible, the latter property is easier to ensure when designing transition probabilities, and most sampling algorithms are indeed reversible with respect to their target distribution.

**Remark 12.4** A simple example of non-reversible Markov chain with invariant probability $Q$ is often obtained in practice by alternating two or more $Q$-reversible transition probabilities. Assume, to simplify, that $\mathcal{B}$ is discrete and that $p_1$ and $p_2$ are transition probabilities that satisfy (12.9). Consider a composite Markov chain for which the transition from $X_n$ to $X_{n+1}$ consists in generating first $Y_n$ according to $p_1(X_n, \cdot)$ and then $X_{n+1}$ according to $p_2(Y_n, \cdot)$. The resulting composite transition probability is

$$p(x,y) = \sum_{z \in \mathcal{B}} p_1(x,z)p_2(z,y).$$

Trivially, $Q$ is invariant by $p$, since it is invariant by $p_1$ and $p_2$, but $p$ is not $Q$-reversible. Indeed, $p$ satisfies (12.7) with

$$\tilde{p}(x,y) = \sum_{z \in \mathcal{B}} p_2(x,z)p_1(z,y).$$

♦

### 12.3.4  Irreducibility and recurrence

While necessary, invariance is not sufficient for a Markov chain to converge to $Q$ in distribution. However, it simplifies the general ergodicity conditions compared to the general theory of Markov chains [147, 160], as summarized below, following [192] (see also [13]). We therefore assume that the transition probability $P$ is such that $Q$ is $P$-invariant.

One says that the Markov chain is $Q$-irreducible (or, simply, irreducible in what follows) if and only if, for all $z \in \mathcal{B}$ and all (measurable) $B \subset \mathcal{B}$ such that $Q(B) > 0$, there exists $n > 0$ with $\mathbb{P}_z(X_n \in B) > 0$. (Irreducibility implies that $Q$ is the only invariant probability of the Markov chain.)

A Markov chain is called periodic if there exists $m > 1$ such that $\mathcal{B}$ can be covered by disjoint subsets $\mathcal{B}_0, \ldots, \mathcal{B}_{m-1}$ that satisfy $P(x, \mathcal{B}_j) = 1$ for all $x \in \mathcal{B}_{j-1}$ if $j \geq 1$ and all $x \in \mathcal{B}_{m-1}$ if $j = 0$. In other terms, the chain loops between the sets $\mathcal{B}_0, \ldots, \mathcal{B}_{m-1}$. If such a decomposition does not exists, the chain is called aperiodic.

A periodic chain cannot satisfy (12.3). Indeed, periodicity implies that $P_x(X_n \in \mathcal{B}_i) = 0$ for all $x \in \mathcal{B}_i$ unless $i = 0 \pmod d$. Since the sets $\mathcal{B}_i$ cover $\mathcal{B}$, (12.3) is only possible with $Q = 0$. Irreducibility and aperiodicity are therefore necessary conditions for ergodicity. Combined with the fact that $Q$ is an invariant probability distribution, these conditions are also sufficient, in the sense that (12.3) is true for $Q$-almost all $x$. (See [192] for a proof.)

Without the knowledge that the chain has an invariant probability, showing convergence usually requires showing that the chain is recurrent, which means that, for any set $B$ such that $Q(B) > 0$, the probability that, starting from $x$, $X_n \in B$ for an infinite number of $n$, written as $\mathbb{P}_x(X_n \in B \ i.o.)$ (for infinitely often) is positive for all $x \in E$ and equal to 1 $Q$-almost surely. The fact that irreducibility and aperiodicity combined with $Q$-invariance imply recurrence (or, more precisely, $Q$-positive recurrent [147]) is an important remark that significantly simplifies the theory for MCMC simulation. Note that, by restricting $\mathcal{B}$ to a suitable set of $Q$-probability 1, one can assume that $\mathbb{P}_x(X_n \in B \ i.o.) = 1$ for all $x \in \mathcal{B}$, which is called *Harris recurrence*. It the chain is Harris recurrent, then (12.3) holds with $\mu_0 = \delta_x$ for all $x \in \mathcal{B}$. [4]

One says that $C \subset \mathcal{B}$ is a "small" set if $Q(C) > 0$ and there exists a triple $(m_0, \epsilon, \nu)$, with $\epsilon > 0$ and $\nu$ a probability distribution on $\mathcal{B}$, such that

$$P^{m_0}(x, \cdot) \geq \epsilon \nu(\cdot)$$

for all $x \in C$. A slightly different result, proved in [13], replaces irreducibility by the property that there exists a small set $C \subset \mathcal{B}$ such that

$$\mathbb{P}_x(\exists n : X_n \in C) > 0$$

---

[4]Harris recurrence is also associated with the uniqueness of right eigenvectors of $P$, that is functions $h : \mathcal{B} \to \mathbb{R}$ such that

$$Ph(x) \overset{\Delta}{=} \int_{\mathcal{B}} P(x, dy)h(y) = h(x).$$

Such functions are also called harmonic for $P$. Because $P$ is a transition probability, constant functions are always harmonic. Harris recurrence, in the current context, is equivalent to the fact that every bounded harmonic function is constant.

for $Q$-almost all $x \in \mathcal{B}$. One then replaces aperiodicity by the similar condition that the greatest common divisor of the set of integers $m$ such that there exists $\epsilon_m$ with $P^m(x, \cdot) \geq \epsilon_m \nu(\cdot)$ for all $x \in C$ is equal to 1.  These two conditions combined with $Q$-invariance also imply that (12.3) holds for $Q$-almost all $x \in \mathcal{B}$.

### 12.3.5   Speed of convergence

It is also important to quantify the speed of convergence in (12.3).  Efficient algorithms typically have a geometric convergence speed, namely

$$D_{\mathrm{var}}(P_x^n, Q) \leq M(x) r^n \tag{12.10}$$

for some $0 \leq r < 1$ and some function $M(x)$, or uniformly geometric convergence speed, for which the function $M$ is bounded (or, equivalently, constant).

A sufficient condition for geometric ergodicity is provided in Nummelin [147, Proposition 5.21].  Assume that the chain is Harris recurrent and that there exist $r > 1$, a small set $C$ and a "drift function" $h$ with

$$\sup_{x \notin C}(r\mathbb{E}(h(X_{n+1}) \mid X_n = x) - h(x)) < 0 \tag{12.11a}$$

and

$$\sup_{x \in C} \mathbb{E}(h(X_{n+1})\mathbf{1}_{X_{n+1} \notin C} \mid X_n = x) < \infty. \tag{12.11b}$$

Then the Markov chain is geometrically ergodic.  Note that $\mathbb{E}(h(X_{n+1}) \mid X_n = x) = Ph(x)$.  Equations (12.11a) and (12.11b) can be summarized in a single equation [136], namely

$$Ph(x) \leq \beta h(x) + M\mathbf{1}_C(x) \tag{12.12}$$

with $\beta = 1/r < 1$ and $M \geq 0$.

### 12.3.6   Models on finite state spaces

Uniform geometric ergodicity is implied by the simple condition that the whole set $\mathcal{B}$ is small, requiring in a uniform lower bound, for some probability distribution $\nu$,

$$P^{m_0}(x, \cdot) \geq \epsilon \nu(\cdot) \tag{12.13}$$

for all $x \in \mathcal{B}$.

Such uniform conditions usually require strong restrictions on the space $\mathcal{B}$, such as compactness or finiteness. To illustrate this consider the case in which the set $\mathcal{B}$ is finite. Assume, to simplify, that $Q(x) > 0$ for all $x \in \mathcal{B}$ (one can restrict the Markov chain to such $x$'s otherwise).  Arbitrarily labeling elements of $\mathcal{B}$ as $\mathcal{B} = \{x_1, \ldots, x_N\}$,

we can consider $p(x,y)$ as the coefficients of a matrix $P = (p(x_k, x_l), k, l = 1, \ldots, N)$. Such a matrix, which has non-negative entries and row sums equal to 1, is called a stochastic matrix.

We will denote the $n$th power of $P$ as $P^n = (p^{(n)}(x_k, x_l), k, l = 1, \ldots, N)$. One immediately sees that irreducibility is equivalent to the fact that, for all $x, y \in \mathcal{B}$, there exists $m$ (that may depend of $x$ and $y$) such that $p^{(m)}(x, y) > 0$. One can furthermore show that the chain is irreducible and aperiodic if one can choose $m$ independent of $x$ and $y$ above, that is, if there exists $m$ such that $P^m$ has positive coefficients. This condition clearly implies uniformly geometric ergodicity, which is therefore valid for all irreducible and aperiodic Markov chains on finite sets.

This result can also be deduced from properties of matrices with non-negative or positive coefficients. The Perron-Frobenius theorem [93] states that the eigenvalue 1 (associated with the eigenvector $\mathbb{1}$) is the largest, in modulus, eigenvalue of a stochastic matrix $\tilde{P}$ with positive entries, that it has multiplicity one and that all other eigenvalues have a modulus strictly smaller that one. If $P^m$ has positive entries, this implies that all the eigenvalues of $(P^m - \mathbb{1}Q)$ (where $Q$ is considered as a row vector) have modulus strictly less than one. This fact can then be used to prove uniformly geometric ergodicity.

### 12.3.7   Examples on $\mathbb{R}^d$

To take a geometrically ergodic example that is not uniform, consider the simple random walk provided by the iterations

$$X_{n+1} = \rho X_n + \tau^2 \epsilon_n$$

where $\epsilon_n \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$, $\tau^2 > 0$ and $0 < \rho < 1$. One shows easily by induction that the conditional distribution of $X_n$ given $X_0 = x$ is Gaussian with mean $m_n = \rho^n x$ and covariance matrix $\sigma_n^2 \mathrm{Id}_{\mathbb{R}^d}$ with

$$\sigma_n^2 = \frac{1 - \rho^{2n}}{1 - \rho^2} \tau^2.$$

In particular, the distribution $Q = \mathcal{N}(0, \sigma_\infty^2 \mathrm{Id}_{\mathbb{R}^d})$, with $\sigma_\infty^2 = \tau^2/(1 - \rho^2)$, is invariant. Estimates on the variational distances between Gaussian distributions, such as those provided in Devroye et al. [61], can then be used to show that

$$D_{\mathrm{var}}(P_x^n, Q) \leq M(x) \rho^n$$

where $M$ grows linearly in $x$ but is not bounded.

Situations in which on can, as above, compute the probability distributions of $X_n$ are rare, however, and proving geometric convergence is significantly more difficult

than for finite-state chains. For chains on $\mathbb{R}^d$ (or, more generally, locally compact metric spaces), the drift function criterion (12.12) can be used. Assume that $Ph(\cdot)$, given by

$$Ph(x) = \mathbb{E}(h(X_{n+1}) \mid X_n = x) = \int_{\mathbb{R}^d} h(y) P(x, dy)$$

is continuous as soon as the function $h : \mathbb{R}^d \to \mathbb{R}$ is continuous (one says that the chain is weak Feller). This true, for example, if $P(x, \cdot)$ has a p.d.f. with respect to Lebesgue's measure which is continuous in $x$. In such a situation, one can see that compact sets are small sets, and (12.12) can be restated as the existence if a positive function $h$ with compact sub-level sets and such that $h(x) \geq 1$, of a compact set $C \subset \mathbb{R}^d$ and of positive constants $\beta < 1$ and $b$ such that, for all $x \in \mathbb{R}^d$,

$$Ph(x) \leq \beta h(x) + b\mathbf{1}_C(x). \tag{12.14}$$

As an example, consider the Markov chain defined by

$$X_{n+1} = X_n - \delta \nabla H(X_n) + \tau \epsilon_{n+1}$$

where $\epsilon_2, \epsilon_2, \ldots$ are i.i.d. standard $d$-dimensional Gaussian variables and $H : \mathbb{R}^d \to \mathbb{R}$ is $C^2$. This chain is clearly irreducible (with respect to Lebesgue's measure). One has

$$Ph(x) = \frac{1}{(2\pi\tau^2)^{d/2}} \int_{\mathbb{R}^d} h(y) e^{-\frac{1}{2\tau^2}|y - x + \delta\nabla H(x)|^2} dy = \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} h(x - \delta\nabla H(x) + \tau u) e^{-\frac{|u|^2}{2}} dy.$$

Let us make the assumption that $H$ is $L$-$C^1$ for some $L > 0$ (c.f. definition 3.15) and furthermore assume that $|\nabla H(x)|$ tends to infinity when $x$ tends to infinity, ensuring the fact that the sets $\{x : |\nabla H(x)| \leq c\}$ are compact for $c > 0$. We want to show that, if $\delta$ is small enough, (12.14) holds for $h(x) = \exp(mH(x))$ and $m$ small enough.

We first compute an upper bound of

$$g(x, u) = mH(x - \delta\nabla H(x) + \tau u) - \frac{|u|^2}{2}.$$

Using the $L$-$C^1$ property, we have

$$g(x, u) \leq mH(x) + m(-\delta\nabla H(x) + \tau u)^T \nabla H(x) + \frac{mL}{2}|\delta\nabla H(x) - \tau u|^2 - \frac{|u|^2}{2}$$

$$= mH(x) - m\delta(1 - \delta L/2)|\nabla H(x)|^2 + m\tau(1 - \tau L)\nabla H(x)^T u - \frac{1 - mL\tau^2}{2}|u|^2$$

$$= mH(x) - \frac{1 - mL\tau^2}{2}\left|u - \frac{m\tau(1 - \tau L)}{1 - mL\tau^2}\nabla H(x)\right|^2$$

$$\quad - m\left(\delta(1 - \delta L/2) - \frac{m\tau^2(1 - \tau L)^2}{2(1 - mL\tau^2)}\right)|\nabla H(x)|^2$$

Assume that $mL\tau^2 \leq 1$. It follows that

$$Ph(x) = \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} e^{g(x,u)} du \leq \frac{h(x)}{(1 - mL\tau^2)^{d/2}}$$

$$\exp\left(-m\left(\delta(1 - \delta L/2) - \frac{m\tau^2(1 - \tau L)^2}{2(1 - mL\tau^2)}\right)|\nabla H(x)|^2\right)$$

Using this upper bound, we see that (12.14) will hold if one first chooses $\delta$ such that $\delta L < 2$, then $m$ such that $mL\tau^2 < 1$ and

$$\frac{m\tau^2(1 - \tau L)^2}{2(1 - mL\tau^2)} < \delta(1 - \delta L/2)$$

and finally choose $C = \{x : |\nabla H(x)| \leq c\}$ where $c$ is large enough so that

$$\frac{1}{(1 - mL\tau^2)^{d/2}} \exp\left(-m\left(\delta(1 - \delta L/2) - \frac{m\tau^2(1 - \tau L)^2}{2(1 - mL\tau^2)}\right)c^2\right) < 1.$$

Note that this Markov chain is not in detailed balance. Since $P(x, \cdot)$ has a p.d.f., being in detailed balance requires the ratio $p(x, y)/p(y, x)$ to simplify as a ratio $q(y)/q(x)$ for some function $q$, which does not hold. However, we can identify the invariant distribution approximately with small $\delta$ and $\tau$, that we will assume to satisfy $\tau = a\sqrt{\delta}$ for a fixed $a > 0$, with $\delta$ a small number.

We can write

$$p(x, y) = \frac{1}{(2\pi\tau^2)^{d/2}} \exp\left(-\frac{1}{2\tau^2}|y - x + \delta\nabla H(x)|^2\right)$$

$$= \frac{1}{(2\pi\tau^2)^{d/2}} \exp\left(-\frac{1}{2\tau^2}|y - x|^2 - \frac{\delta}{\tau^2}(y - x)^T\nabla H(x) - \frac{\delta^2}{2\tau^2}|\nabla H(x)|^2\right).$$

If $q$ is a density, we have

$$qP(y) = \int_{\mathbb{R}^d} q(x)p(x, y)dx$$

$$= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} q(y + a\sqrt{\delta}u)\exp\left(-\frac{1}{2}|u|^2 + \frac{\sqrt{\delta}}{a}u^T\nabla H(y + a\sqrt{\delta}u) - \frac{\delta}{2a^2}|\nabla H(y + a\sqrt{\delta}u)|^2\right)du$$

Make the expansions:

$$q(y + a\sqrt{\delta}u) = q(y) + a\sqrt{\delta}\nabla q(y)^T u + \frac{a^2\delta}{2}u^T\nabla^2 q(y)u + o(\delta|u|^2)$$

and

$$\exp\left(\frac{\sqrt{\delta}}{a}u^T\nabla H(y+a\sqrt{\delta}u) - \frac{\delta}{2a^2}|\nabla H(y+a\sqrt{\delta}u)|^2\right)$$

$$= 1 + \frac{\sqrt{\delta}}{a}u^T\nabla H(y) - \frac{\delta}{2a^2}|\nabla H(y)|^2 + \delta u^T\nabla^2 H(u)u + \frac{\delta}{2a^2}(u^T\nabla H(y))^2 + o(\delta|u|^2).$$

Taking the product and using the fact that $(2\pi)^{-d/2}\int_{\mathbb{R}^d}u\exp(-|u|^2/2)du = 0$ and that $(2\pi)^{-d/2}\int_{\mathbb{R}^d}u^T Au\exp(-|u|^2/2)du = \text{trace}(A)$ for any symmetric matrix $A$, we can write, taking the product:

$$qP(y) = q(y) + \delta\left(\frac{a^2}{2}\Delta q(y) + \nabla H(y)^T\nabla q(y) + q(y)\Delta H(y)\right) + o(\delta)$$

This indicates that, if $q$ is invariant by $P$, it should satisfy

$$\frac{a^2}{2}\Delta q(y) + \nabla H(y)^T\nabla q(y) + q(y)\Delta H(y) = o(1).$$

The partial differential equation

$$\frac{a^2}{2}\Delta q(y) + \nabla H(y)^T\nabla q(y) + q(y)\Delta H(y) = 0 \tag{12.15}$$

is satisfied by the function $y \mapsto e^{-\frac{2H(y)}{a^2}}$. Assuming that this function is integrable, this computation suggests that, for small $\delta$, the Markov chain approximately samples from the probability distribution

$$q_0 = \frac{1}{Z}e^{-\frac{2H(x)}{a^2}}.$$

This is further discussed in the next remark that involves stochastic differential equations. We will also present a correction of this Markov chain that samples from $q_0$ for all $\delta$ in section 12.5.2.

**Remark 12.5 (Langevin equation)**  This chain is indeed the Euler discretization [106] of the stochastic differential equation,

$$dx_t = -\nabla H(x_t)dt + adw_t \tag{12.16}$$

where $w_t$ is a Brownian motion. Under general hypotheses, this stochastic diffusion equation, called a *Langevin equation*, indeed converges in distribution to $q_0(x)$. [5]

---

[5]Providing a rigorous account of the theory of stochastic differential equations is beyond our scope, and we refer the reader to the many textbooks on the subject, such as McKean [131], Ikeda and Watanabe [96], Ethier and Kurtz [69] (see also Berglund [26] for a short introduction).

Such diffusions are continuous-time Markov processes $(X_t, t \geq 0)$, which means that the probability distribution of $X_{t+s}$ given all events before and including time $s$ only depends on $X_s$ and is provided by a transition probability $P_t$, with

$$\mathbb{P}(X_{t+1} \in A \mid X_s = x) = P_t(x, A).$$

Similarly to deterministic ordinary differential equations, one shows that under sufficient regularity conditions (e.g., $\nabla H$ is $C^1$), equations such as (12.16) have solutions up to some positive (random) explosion time, and that this explosion time is finite under additional conditions that ensure that $|\nabla H(x)|$ does not grow too fast when $x$ tends to infinity.

If $\varphi$ is a smooth enough function (say, $C^2$, with compact support), the function $(t, x) \mapsto P_t \varphi(x)$ satisfies the partial differential equation, called Kolmogorov's backward equation,

$$\partial_t P_t \varphi(x) = -\nabla H(x)^T \nabla P_t \varphi(x) + \frac{a^2}{2} \Delta P_t \varphi(x)$$

with initial condition $P_0 \varphi(x) = \varphi(x)$. If $P_t(x, \cdot)$ has at all times $t$ a p.d.f. $p_t(x, \cdot)$, then this p.d.f. must satisfy the forward Kolmogorov equation:

$$\partial_t p_t(x, y) = \nabla_2 \cdot (\nabla H(y) p_t(x, y)) + \frac{a^2}{2} \Delta_2 p_t(x, y)$$

where $\nabla_2$ and $\Delta_2$ indicate differentiation with respect to the second variable $(y)$. (Recall that $\delta f$ denotes the Laplacian of $f$.) Moreover, if $Q$ is an invariant distribution with p.d.f. $q$, it satisfies the equation

$$\nabla \cdot (q \nabla H) + \frac{a^2}{2} \Delta q(y) = 0.$$

Noting that $\nabla \cdot (q \nabla H) = \nabla q^T \nabla H + q \Delta H$, we retrieve (12.15). Convergence properties (and, in particular, geometric convergence) of the Langevin equation to its limit distribution are studied in Roberts and Tweedie [166], using methods introduced in Meyn and Tweedie [134, 135, 136] ⧫

## 12.4 Gibbs sampling

### 12.4.1 Definition

The Gibbs sampling algorithm [79] was introduced to sample from a distribution on large sets for which direct sampling is intractable and rejection samping is inefficient. It generates a Markov chain that converges (under some hypotheses) in distribution to this target probability. A general version of this algorithm is described below.

Let $Q$ be a probability distribution on $\mathcal{B}$. Consider a finite family $U_1, \ldots, U_K$ of random variables defined on $\mathcal{B}$ with values in measurable spaces $\mathcal{B}'_1, \ldots, \mathcal{B}'_K$. Let $Q^i = Q_{U^i}$ denote the image of $Q$ by $U^i$, defined by $Q^i(B_i) = Q(U_i \in B_i)$ for $B_i \subset \mathcal{B}_i$. Also, assume that there exists, for all $i$, a regular family of conditional probabilities for $Q$ given $U_i$, defined as a collection of transition probabilities $(u_i, A) \mapsto Q_i(u_i, A)$ for $u_i \in \mathcal{B}_i$ and $A \subset \mathcal{B}$, that satisfy

$$\int_A g(U_i(x)) Q(dx) = \int_{\mathcal{B}_i} Q_i(u_i, A) g(u_i) Q^i(du_i)$$

for all nonnegative measurable functions $g : \mathcal{B}_i \to \mathbb{R}$. In simpler terms, $Q_i(u_i, A)$ determine a consistent set of conditional probabilities for $Q(A \mid U_i = u_i)$. For discrete random variables (resp. variables with p.d.f.'s on $\mathbb{R}^d$), they are just elementary conditional probabilities.

We then consider the following algorithm.

---

**Algorithm 12.2 (Gibbs sampling)**
Initialize the algorithm with some $z(0) = z_0 \in \mathcal{B}$ and iterate the following two update steps given a current $z(n) \in \mathcal{B}$:

(1)  Select $j \in \{1, \ldots, K\}$ according to some pre-defined scheme, i.e., at random according to a probability distribution $\pi^{(n)}$ on the set $\{1, \ldots, K\}$.

(2)  Sample a new value $z(n+1)$ according to the probability distribution $Q_j(U_j(z(n)), \cdot)$.

---

One typically chooses the probability distribution in Step 1 equal to the uniform distribution on $\{1, \ldots, K\}$ (in which case it is independent on $n$), or to $\pi^{(n)} = \delta_{j_n}$ where $j_n = 1 + (n \,(mod)\, K)$ (periodic scheme). Strictly speaking, Gibbs sampling is a Markov chain if $\pi^{(n)}$ does not depend on $n$, and we will make this simplifying assumption in the rest of our discussion (therefore replacing $\pi^{(n)}$ by $\pi$). One obvious requirement for the feasibility of the method is that step (2) can be performed efficiently since it must be repeated a very large number of times.

One can see that the Markov chain generated by this algorithm is $Q$-reversible. Indeed, assume that $X_n \sim Q$. For any (measurable) subsets $A$ and $B$ in $\mathcal{B}$, one has, using the definition of conditional expectations,

$$\mathbb{P}(X_n \in A, X_{n+1} \in B) = \sum_{i=1}^{K} \mathbb{E}\left(\mathbf{1}_{Z \in A} Q_i(U_i(Z), B)\right) \pi(i). \tag{12.17}$$

Now, for any $i$

$$\mathbb{E}\Big(\mathbf{1}_{Z\in A}Q_i(U_i(Z),B)\Big) = \int_A Q_i(U_i(z),B)Q(dz)$$

$$= \int_{\mathcal{B}_i} Q_i(u_i,A)Q_i(u_i,B)Q^i(du_i)$$

which is symmetric in $A$ and $B$.

Note that, in the discrete case

$$P(z,\tilde{z}) = \sum_{i=1}^n \pi(i)\frac{Q(\tilde{z})\mathbf{1}_{U^i(\tilde{z})=U^i(z)}}{\sum_{z':U^i(z')=U^i(z)} Q(z')} \tag{12.18}$$

and the relation $Q(z)P(z,\tilde{z}) = Q(\tilde{z})P(\tilde{z},z)$ is obvious.

The conditioning variables $U_1,\dots,U_K$ should ensure, at least, that the associated Markov chain is irreducible and aperiodic. For irreducibility, this requires that $Z$ can visits $Q$-almost all elements of $\mathcal{B}$ by a sequence of steps that lead one of the $U_i$'s invariant.

**Remark 12.6** In the standard version of Gibbs sampling, $\mathcal{B}$ is a product space $\mathcal{B}_1 \times \cdots \times \mathcal{B}_K$, and

$$\mathcal{B}'_j = \mathcal{B}_1 \times \cdots \times \mathcal{B}_{j-1} \times \mathcal{B}_{j+1} \times \cdots \times \mathcal{B}_K.$$

One then takes $U_j(z^{(1)},\dots,z^{(K)}) = (z^{(1)},\dots,z^{(i-1)},z^{(i+1)},\dots,z^{(K)})$. In other terms, step 2 in the algorithm replaces the current value of $z^{(j)}(n)$ by a new one sampled from the conditional distribution of $Z^{(j)}$ given the current values of $z^{(i)}(n), i \neq j$. ♦

**Remark 12.7** We have considered a fixed number of conditioning variables, $U_1,\dots,U_K$, for simplicity, but the same analysis can be carried on if one replaces $U_j$ by a function $U : (x,\theta) \mapsto U_\theta(x)$ defined on a product space $\mathcal{B} \times \Theta$, taking values in some space $\tilde{\mathcal{B}}$, where $\Theta$ is a probability space equipped with a probability distribution $\pi$ and $\mathcal{B}$ is measurable. The previous discussion corresponds to $\Theta = \{1,\dots,K\}$ and $\mathcal{B} = \bigcup_{i=1}^K \{i\} \times \mathcal{B}_i$ (so that $U_i(x)$ is replaced by $(i,U_i(x))$).

One may then define $Q^\theta$ as the image of $Q$ by $U_\theta$ and let $Q_\theta(u,A)$ provide a version of $Q(A \mid U_\theta = u)$. The only change in the previous discussion (besides using $\theta$ in index) is that (12.17) becomes

$$\mathbb{P}(X_n \in A, X_{n+1} \in B) = \int_\Theta \mathbb{E}\Big(\mathbf{1}_{Z\in A}Q_\theta(U_\theta(Z),B)\Big)\pi(d\theta).$$

♦

**Remark 12.8** Using notation from the previous remark, and allowing $\pi = \pi^{(n)}$ to depend on $n$, it is possible to allow $\pi^{(n)}$ to depend on the current state $z(n)$ using the following construction.

For every step $n$, assume that there exists a subset $\Theta_n$ of $\Theta$ such that $\pi^{(n)}(z, \Theta_n) = 1$ and that, for all $\theta \in \Theta_n$, $\pi^{(n)}$ can be expressed in the form

$$\pi^{(n)}(z, \cdot) = \psi_\theta^{(n)}(U_\theta(z), \cdot)$$

for some transition probability $\psi_\theta^{(n)}$ from $\mathcal{B}_\theta$ to $\Theta_n$. The resulting chain remains $Q$-reversible, since

$$\mathbb{P}(X_n \in A, X_{n+1} \in B) = \int_{\mathcal{B}} \int_{\Theta_n} \mathbf{1}_{z \in A} Q_\theta(U_\theta(z), B) \pi^{(n)}(z, d\theta) Q(dz)$$

$$= \int_{\Theta_n} \int_{\mathcal{B}} \mathbf{1}_{z \in A} Q_\theta(U_\theta(z), B) \psi_\theta^{(n)}(U_\theta(z), d\theta) Q(dz)$$

$$= \int_{\Theta_n} \int_{\tilde{\mathcal{B}}} Q_\theta(u, A) Q_\theta(u, B) \psi_\theta^{(n)}(u, d\theta) Q^\theta(du). \qquad \blacklozenge$$

### 12.4.2   Example: Ising model

We will see several examples of applications of Gibbs sampling in the next few chapters. Here, we consider a special instance of Markov random field (see chapter 13) called the Ising model. For this example, $\mathcal{B} = \{0, 1\}^L$, and

$$q(z) = \frac{1}{C} \exp\left( \sum_{j=1}^{L} \alpha z^{(j)} + \sum_{i,j=1, i<j}^{L} \beta_{ij} z^{(i)} z^{(j)} \right).$$

Note that, although $\mathcal{B}$ is a finite set, its cardinality, $2^L$, is too large for the enumerative procedure described in section 12.1 to be applicable as soon as $L$ is, say, larger than 30. In practical applications of this model, $L$ is orders of magnitude larger, typically in the thousands or tens of thousands.

We here apply standard Gibbs sampling, as described in remark 12.6, defining $\mathcal{B}_j = \{0, 1\}$ and
$$U_i(z^{(1)}, \ldots, z^{(L)}) = (z^{(1)}, \ldots, z^{(i-1)}, z^{(i+1)}, \ldots, z^{(L)}).$$

The conditional distribution of $Z^{(j)}$ given $U_j(z)$ is a Bernoulli distribution with parameter
$$q_{Z^{(j)}}(1 \mid U_j(z)) = \frac{\exp(\alpha + \sum_{j'=1, j' \neq j}^{L} \beta_{jj'} z^{(j')})}{1 + \exp(\alpha + \sum_{j'=1, j' \neq j}^{L} \beta_{ij} z^{(j)})}$$

(taking $\beta_{jj'} = \beta_{j'j}$ for $j > j'$). Gibbs sampling for this model will generate a sequence of variables $Z(0), Z(1), \ldots$ by fixing $Z(0)$ arbitrarily and, given $Z(n) = z$, applying the two steps:

1. Select $j \in \{1, \ldots, L\}$ at random according to a probability distribution $\pi^{(n)}$ on the set $\{1, \ldots, L\}$.

2. Sample a new value $\zeta \in \{0, 1\}$ according to the Bernoulli distribution with parameter $q_{Z_n^{(j)}}(1 \mid U_j(z))$, and set $Z^{(j)}(n+1) = \zeta$ and $Z^{(j')}(n+1) = Z^{(j')}(n)$ for $j' \neq j$.

Let us now consider the Ising model with fixed total activation, namely the previous distribution conditional to $S(z) \stackrel{\Delta}{=} z^{(1)} + \cdots + z^{(L)} = h$ where $0 < h < L$. The distribution one wants to sample from now is

$$q_h(z) = \frac{1}{C_h} \exp\left( \sum_{j=1}^{L} \alpha z^{(j)} + \sum_{i,j=1, i<j}^{L} \beta_{ij} z^{(i)} z^{(j)} \right) \mathbf{1}_{S(z)=h}.$$

In that case, the previous choice for the one-step transitions does not work, because fixing all but one coordinate of $z$ also fixes the last one (so that the chain would not move from its initial value and would certainly not be irreducible). One can however fix all but two coordinates, therefore defining

$$U_{ij}(z^{(1)}, \ldots, z^{(L)}) = (z^{(1)}, \ldots, z^{(i-1)}, z^{(i+1)}, \ldots, z^{(j-1)}, z^{(j+1)}, \ldots, z^{(L)})$$

and $B_{ij} = \{0, 1\}^2$. If $U_{ij}(z)$ is fixed, the only acceptable configurations are $z$ and the configuration $z'$ deduced from $z$ by switching the value of $z^{(i)}$ and $z^{(j)}$. Thus, there is no possible change is $z^{(i)} = z^{(j)}$. If $z^{(i)} \neq z^{(j)}$, then the probability of flipping the values of $z^{(i)}$ and $z^{(j)}$ is $q_h(z')/(q_h(z) + q_h(z'))$.

## 12.5 Metropolis-Hastings

### 12.5.1 Definition

Gibbs sampling is a special case of a generic MCMC algorithm called Metropolis-Hastings that is defined as follows [133, 88]. Assume that the distribution $Q$ has a density $q$ with respect to a measure $\mu$ on $\mathcal{B}$. Specify a transition probability on $\mathcal{B}$, represented by a family of density functions with respect to $\mu$, $(g(z, \cdot), z \in \mathcal{B})$, and a family of acceptance functions $(z, z') \mapsto a(z, z') \in [0, 1]$. Two basic examples are when $\mathcal{B}$ is finite, $\mu$ is the counting measure, and $q$ and $g$ are probability mass functions, and when $\mathcal{B} = \mathbb{R}^d$, $\mu$ is Lebesgue's measure and $q$ and $g$ are probability density functions.

The sampling algorithm is then defined as follows. It invokes a function $a$ that will be specified below.

**Algorithm 12.3 (Metropolis-Hastings)**
Initialize the algorithm with $Z(0) = z(0) \in \mathcal{B}$. At step $n$, the current value $Z(n) = z$ is then updated as follows.

- "Propose" a new configuration $z'$ drawn according to $g(z, \cdot)$.

- "Accept" $z'$ (i.e., set $Z(n + 1) = z'$) with probability $a(z, z')$. If the new value is rejected, keep the current one, i.e., let $Z(n + 1) = z$.

---

The transition probabilities for this process are $p(x, y) = g(x, y)a(x, y)$ if $x \neq y$ and $p(x, x) = 1 - \sum_{y \neq x} p(x, y)$. The chain is $Q$-reversible is the detailed balance equation

$$q(z)g(z, z')a(z, z') = q(z')g(z', z)a(z', z) \tag{12.19}$$

is satisfied. The functions $g$ and $a$ are part of the design of the algorithm, but (12.19) suggest that $g$ should satisfy the "weak symmetry" condition:

$$\forall x, y \in \Omega : g(x, y) = 0 \Leftrightarrow g(y, x) = 0. \tag{12.20}$$

Note that this condition is necessary to ensure (12.19) if $q(z) > 0$ for all $z$. If $q(z) > 0$, the fact that acceptance probabilities are less than 1 requires that

$$a(z, z') \leq \min\left(1, \frac{q(z')g(z', z)}{q(z)g(z, z')}\right).$$

If one takes $a(z, z')$ equal to the r.h.s., so that

$$a(z, z') = \min\left(1, \frac{q(z')g(z', z)}{q(z)g(z, z')}\right), \tag{12.21}$$

then (12.19) is satisfied as soon as $q(z) > 0$. If $q(z) = 0$, then this definition ensures that $a(z', z) = 0$ and (12.19) is also satisfied. Note also that the case $g(z, z') = 0$ is not relevant, since $z'$ is not attainable from $z$ in one step in this case. This shows that (12.21) provides a $Q$-reversible chain. Obviously, if $g$ already satisfies $q(z)g(z, z') = q(z')g(z', z)$, which is the case for Gibbs sampling, then one should take $a(z, z') = 1$ for all $z$ and $z'$.

### 12.5.2   Sampling methods for continuous variables

While the Gibbs sampling and Metropolis-Hastings methods can be (and were) formulated for general variables and probability distributions, proving that the related chains are ergodic, and checking conditions for geometric convergence speed is much harder when dealing with general state spaces than with finite or compact spaces (see, e.g., [164, 132, 6, 165]). On the other hand, interesting choices

of proposal transitions for Metropolis-Hastings are available when $\mathcal{B} = \mathbb{R}^d$ and $\mu$ is Lebesgue's measure, taking advantage, in particular, of differential calculus. More precisely, assume that $q$ takes the form

$$q(z) = \frac{1}{C} \exp(-H(z))$$

for some smooth function $H$ (at least $C^1$), such that $\exp(-H)$ is integrable. We saw in section 12.3.7 that, under suitable assumptions, the Markov chain

$$X_{n+1} = X_n - \frac{\delta}{2} \nabla H(X_n) + \sqrt{\delta} \epsilon_{n+1} \tag{12.22}$$

with $\epsilon_{n+1} \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$ has $q$ as invariant distribution in the limit $\delta \to 0$. Its transition probability, such that $g(z, \cdot)$ is the p.d.f. of $\mathcal{N}(z - \frac{\delta}{2} \nabla H(z), \delta \mathrm{Id}_{\mathbb{R}^d})$, is therefore a natural choice for a proposal distribution in the Metropolis-Hastings algorithm. In addition to converging from the exact target distribution, this "Metropolis Adjusted Langevin Algorithm" (or MALA) can also be proved to satisfy geometric convergence under less restrictive hypotheses than (12.22) [166].

Another approach, similar to MALA is the Hamiltonian Monte-Carlo methods (or hybrid Monte-Carlo) [65, 142]. Inspired by physics, the method introduces a new variable, $p \in \mathbb{R}^d$, called "momentum," and defines the "Hamiltonian:"

$$\mathcal{H}(z, m) = H(z) + \frac{1}{2}|m|^2.$$

Fix a time $\theta > 0$. The proposal transition $g(z, \cdot)$ is then defined as the value $\zeta(\theta)$ that is obtained by solving the Hamiltonian dynamical system

$$\begin{cases} \partial_t \zeta(t) = \partial_p \mathcal{H}(\zeta(t), \mu(t)) = \mu(t) \\ \partial_t \mu(t) = -\partial_z \mathcal{H}(\zeta(t), \mu(t)) = -\nabla H(\zeta(t)) \end{cases} \tag{12.23}$$

with $\zeta(0) = z$ and $\mu(0) \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$. One can easily see that $\partial_t \mathcal{H}(\zeta(t), \mu(t)) = 0$, which implies that

$$H(\zeta(t)) + \frac{1}{2}|\mu(t)|^2 = H(z) + \frac{1}{2}|\mu(0)|^2$$

at all times $t$, or, denoting by $\varphi_{\mathcal{N}}$ the p.d.f. of the $d$-dimensional standard Gaussian,

$$q(\zeta(t))\varphi_{\mathcal{N}}(\mu(t)) = q(\zeta(0))\varphi_{\mathcal{N}}(\mu(0)).$$

Moreover, if one denotes by $\Phi_t(z, m) = (\boldsymbol{z}_t(z, m), \boldsymbol{m}_t(z, m))$ the solution $(\zeta(t), \mu(t))$ of the system started with $\zeta(0) = z$ and $\mu(0) = m$, one can also see that $\det(d\Phi_t(z, m)) = 1$ at all times. Indeed, applying (1.5) and the chain rule, we have

$$\partial_t \log \det(d\Phi_t(z, m)) = \mathrm{trace}(d\Phi_t(z, m)^{-1} \partial_t d\Phi_t(t, m)).$$

From

$$\begin{cases} \partial_t \boldsymbol{z}_t(z,m) = \boldsymbol{m}_t(z,p) \\ \partial_t \boldsymbol{p}_t(z,m) = -\nabla H(\boldsymbol{z}_t(z,m)) \end{cases}$$

we get

$$\partial_t d\Phi_t(z,m) = \begin{pmatrix} \partial_z \boldsymbol{m}_t(z,m) & \partial_m \boldsymbol{m}_t(z,m) \\ -\nabla^2 H(\boldsymbol{z}_t(z,m))\partial_z \boldsymbol{z}_t(z,m) & -\nabla^2 H(\boldsymbol{z}_t(z,m))\partial_m \boldsymbol{z}_t(z,m) \end{pmatrix}$$

$$= \begin{pmatrix} 0 & \mathrm{Id}_{\mathbb{R}^d} \\ -\nabla^2 H(\boldsymbol{z}_t(z,m)) & 0 \end{pmatrix} d\Phi_t(z,m).$$

We therefore get

$$\partial_t \log \det(d\Phi_t(z,m)) = \mathrm{trace}\begin{pmatrix} 0 & \mathrm{Id}_{\mathbb{R}^d} \\ -\nabla^2 H(\boldsymbol{z}_t(z,m)) & 0 \end{pmatrix} = 0$$

showing that the determinant is constant. Since $\Phi_0(z,m) = (z,m)$ by definition, we get $\det(d\Phi_t(z,m)) = 1$ at all times.

Let $\bar{q}_t$ denote the p.d.f. of $\Phi_t(z,m)$ and assume that $\bar{q}_0(z,m) = q(z)\varphi_{\mathcal{N}}(m)$. We have, using the change of variable formula

$$\bar{q}_t(\Phi_t(z,m))|\det d\Phi_t(z,m)| = q(z)\varphi_{\mathcal{N}}(m)$$

but the r.h.s. is, from the remarks above also equal to

$$q(\boldsymbol{z}_t(z,m))\varphi_{\mathcal{N}}(\boldsymbol{m}_t(z,m))|\det d\Phi_t(z,m)|$$

yielding the identification

$$\bar{q}_t(z',m') = q(z')\varphi_{\mathcal{N}}(m')$$

This shows that $Q$ (with p.d.f. $q$) is left invariant by this Markov chain. One can actually show that chain is in detailed balance for the joint density $\bar{q}(z,m) = q(z)\varphi_{\mathcal{N}}(m)$. This is due to the fact that the system (12.23) is reversible, in the sense that

$$\Phi_t(\boldsymbol{z}_t(z,m), -\boldsymbol{m}_t(z,m)) = (z,-m),$$

i.e., the system solved from its end point after changing the sign of the momentum returns to its initial state after changing the sign of the momentum a second time. In other terms, letting $J(z,m) = (z,-m)$, we have $\Phi_t^{-1} = J\Phi_t \circ J$. So, consider a function $f : (\mathbb{R}^d \times \mathbb{R}^d)^2 \to \mathbb{R}$. Denoting the Markov chain by $(Z_n, M_n)$, we assume that the next pair $Z_{n+1}, M_{n+1}$ is computed by (i) sampling $M'_n \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$; (ii) solving (12.23), with initial conditions $\zeta(0) = Z_n$ and $\mu(0) = M'_n$; (iii) taking $Z_{n+1} = \zeta(\theta)$ and sampling $M_{n+1} \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$.

We have

$$\mathbb{E}(f(Z_n, M_n, Z_{n+1}, M_{n+1})) = \int f(z, \tilde{m}, \boldsymbol{z}(z,m), \bar{m})\varphi_{\mathcal{N}}(m)\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(z)\,dm\,d\bar{m}\,d\tilde{m}\,dz.$$

Make the change of variables $z' = \mathbf{z}(z, m)$, $m' = \mathbf{m}(z, m)$, which has Jacobian determinant 1, and is such that $z = \mathbf{z}(z', -m')$, $m = -\mathbf{m}(z', -m')$. We get

$$\mathbb{E}(f(Z_n, M_n, Z_{n+1}, M_{n+1}))$$
$$= \int f(\mathbf{z}(z', -m'), \tilde{m}, z', \bar{m})\varphi_{\mathcal{N}}(-\mathbf{m}(z', -m'))\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(\mathbf{z}(z', -m'))dm'd\bar{m}d\tilde{m}dz'$$
$$= \int f(\mathbf{z}(z', -m'), \tilde{m}, z', \bar{m})\varphi_{\mathcal{N}}(\mathbf{m}(z', -m'))\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(\mathbf{z}(z', -m'))dm'd\bar{m}d\tilde{m}dz'$$
$$= \int f(\mathbf{z}(z', -m'), \tilde{m}, z', \bar{m})\varphi_{\mathcal{N}}(-m')\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(z')dm'd\bar{m}d\tilde{m}dz',$$

using the conservation of $\mathcal{H}$. Making the change of variables $m' \to -m'$, we get

$$\mathbb{E}(f(Z_n, M_n, Z_{n+1}, M_{n+1})$$
$$= \int f(\mathbf{z}(z', m'), \tilde{m}, z', \bar{m})\varphi_{\mathcal{N}}(m')\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(z')dm'd\bar{m}d\tilde{m}dz'$$

which is equal to $\mathbb{E}(f(Z_{n+1}, M_{n+1}, Z_n, M_n))$ showing the reversibility of the chain.

This simulation scheme can potentially make large moves in the current configuration $z$ while maintaining detailed balance (therefore not requiring an accept/reject step). However, practical implementations require discretizing (12.23), which breaks the conservation properties that were used in the argument above, therefore requiring a Metropolis-Hastings correction. For example, a second-order Runge Kutta (RK2) scheme with time step $\alpha$ gives

$$\begin{cases} Z_{n+1} = Z_n + \alpha M_n - \dfrac{\alpha^2}{2}\nabla H(Z_n) \\ M_{n+1} = M_n - \dfrac{\alpha}{2}(\nabla H(Z_n) + \nabla H(Z_n + hM_n)) \end{cases}$$

Only the update for $Z_n$ matters, however, since $M_{n+1}$ is discarded and resampled at each step. Importantly, if we let $\delta = \sqrt{\alpha}$ the first equation in the system becomes

$$Z_{n+1} = Z_n - \frac{\delta}{2}\nabla H(Z_n) + \delta M_n$$

with $M_n \sim \mathcal{N}(0, 1)$, which is exactly (12.22). Note that one can, in principle, solve (12.23) for more that one discretization step (the continuous equation can be solved for an arbitrary time), but one must then face the challenge of computing the Metropolis correction since the Hamiltonian is not conserved at each step.

One can however use schemes that are more adapted to solving Hamiltonian

systems [119], such as the Störmer-Verlet scheme, which is

$$
\begin{cases}
M_{n+1/2} = M_n - \dfrac{\alpha}{2}\nabla H(Z_n) \\
\quad Z_{n+1} = Z_n + \alpha M_{n+1/2} \\
\quad M_{n+1} = M_{n+1/2} - \dfrac{\alpha}{2}\nabla H(Z_{n+1})
\end{cases}
$$

This scheme computes $\psi_1 \circ \psi_2 \circ \psi_1(z,m)$ with $\psi_1(z,m) = (z, m - (\alpha/2)\nabla H(z))$ and $\psi_2(z,m) = (z + \alpha m, m)$. Because both $\psi_1$ and $\psi_2$ have a Jacobian determinant equal to 1, so does their composition. This scheme is also reversible, since we have

$$
\begin{cases}
-M_{n+1/2} = -M_{n+1} - \dfrac{\alpha}{2}\nabla H(Z_{n+1}) \\
\quad Z_n = Z_{n+1} - \alpha M_{n+1/2} \\
\quad -M_n = -M_{n+1/2} - \dfrac{\alpha}{2}\nabla H(Z_n)
\end{cases}
$$

These properties are conserved if one applies the Störmer-Verlet scheme more than once at each iteration, that is, fixing some $N > 0$ and letting $\Phi(z,m) = (\psi_1 \circ \psi_2 \circ \psi_1)^{\circ N}$, then $\Phi^{-1} = J\Phi \circ J$, with $J(z,m) = (z,-m)$ with $\det d\Phi = 1$. Considering again the augmented chain which, starting from $(Z_n, M_n)$, samples $\tilde{M} \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$, then computes $(Z', \tilde{M}') = \Phi(Z_n, \tilde{M})$ and finally samples $M' \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$ as a Metropolis-Hastings proposal to sample from $(z,m) \mapsto q(z)\varphi_{\mathcal{N}}(m)$, then, assuming that $(Z,M)$ follows this target distribution and letting $(Z', M')$ be the result of the proposal distribution, we have, as computed above

$$
\mathbb{E}(f(Z,M,Z',M'))
$$
$$
= \int f(z,\tilde{m},\mathbf{z}(z,m),\bar{m})\varphi_{\mathcal{N}}(m)\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(z)\,dm\,d\bar{m}\,d\tilde{m}\,dz
$$
$$
= \int f(\mathbf{z}(z',m'),\tilde{m},z',\bar{m})\varphi_{\mathcal{N}}(\mathbf{m}(z',m'))\varphi_{\mathcal{N}}(\bar{m})\varphi_{\mathcal{N}}(\tilde{m})q(\mathbf{z}(z',m'))\,dm'\,d\bar{m}\,d\tilde{m}\,dz'
$$

This shows that the acceptance probability in the Metropolis step is

$$
a(z,m,z',m') = \min\left(1, \frac{\varphi_{\mathcal{N}}(\mathbf{m}(z',m'))q(\mathbf{z}(z',m'))}{\varphi_{\mathcal{N}}(m)q(z)}\right)
$$
$$
= \exp\left(-\max\left(\mathcal{H}(\mathbf{z}(z',m'),\mathbf{m}(z',m')) - \mathcal{H}(z,m)\right), 0\right)
$$

While the Hamiltonian is not kept invariant by the Störmer-Verlet scheme, so that an accept-reject step is needed, it is usually quite stable over extended periods of time so that the acceptance probability is generally close to one.

## 12.6   Perfect sampling methods

The Markov chain simulation methods, provided in the previous sections do not provide exact samples from the distribution $q$, but only increasingly accurate ap-

proximations. Perfect sampling algorithms [156, 157, 71] use Markov chains "backwards" to generate exact samples. To describe them, it is easier to describe a Markov chain as a stochastic recursive equation of the form

$$X_{n+1} = f(X_n, U_{n+1}) \tag{12.24}$$

where $U_{n+1}$ is independent of $X_n, X_{n-1}, \ldots$, and the $U_k$'s are identically distributed. In the discrete case (assumed in this section), and given a stochastic matrix $P$, one can take $U_n$ to be the uniformly distributed variable used to sample from $(p(X_n, x), x \in \mathcal{B})$. Conversely, the transition probability associated to (12.24) is $p(x, y) = P(f(x, U) = y)$.

It will be convenient to consider negative times also. For $n > 0$, recursively define $F_{-n}(x, u_{-n+1}, \ldots, u_0)$ by

$$F_{-n-1}(x, u_{-n}, \ldots, u_0) = F_{-n}(f(x, u_{-n}), u_{-n+1}, \ldots, u_0)$$

and $F_{-1}(x, u_0) = f(x, u_0)$. Denote, for short, $U_{-n}^0 = (U_{-n}, \ldots, U_0)$. The function $F_{-n}(x, u_{-n+1}^0)$ provides the value of $X_0$ when $X_{-n} = x$ and $U_{-n+1}^0 = u_{-n+1}^0$.

For an infinite past sequence, $u_{-\infty}^0$, let $\nu(u_{-\infty}^0)$ denote the first integer $n$ such that $F_{-n}(x, u_{-n+1}^0)$ does not depend on $x$ (the function "coalesces"). Then, the following theorem is true:

**Theorem 12.9** *Assume that the chain defined by* (12.24) *is ergodic, with invariant distribution Q. Then* $\nu = \nu(U_{-\infty}^0)$ *is finite with probability 1, and*

$$X_* := F_{-\nu}(x, U_{-\nu+1}^0) \tag{12.25}$$

*(which is independent of x) has distribution Q.*

Proof Because the chain is ergodic, we know that there exists an integer $N$ such that one can pass from any state to any other with positive probability. So the chain can, starting from anywhere, coalesce with positive probability in $N$ steps; $\nu$ being infinite would imply that this event never occurs in an infinite number of trials, and this has probability 0.

For any $k > 0$ and any $x \in \mathcal{B}$, we have

$$X_* = F_{-\nu}(f_{-k}(x, U_{-\nu-k+1}^{-\nu}), U_{-\nu+1}^0) = F_{-\nu-k}(x, U_{-\nu-k+1}^0). \tag{12.26}$$

But, because the chain is ergodic, we have, for any $x \in \mathcal{B}$

$$\lim_{k \to \infty} \mathbb{P}(F_{-k}(x, U_{-k+1}^0) = y) = Q(y).$$

We can write

$$\mathbb{P}(F_{-k}(x, U_{-k+1}^0) = y) = \mathbb{P}(F_{-k}(x, U_{-k+1}^0) = y, \nu \le k) + \mathbb{P}(F_{-k}(x, U_{-k+1}^0) = y, \nu > k)$$

$$= \mathbb{P}(X_* = y, \nu \le k) + \mathbb{P}(F_{-k}(x, U_{-k+1}^0) = y, \nu > k)$$

The right-hand side tends to $\mathbb{P}(X_* = y)$ when $k$ tends to infinity (because $\mathbb{P}(v > k)$ tends to 0), and the left-hand side tends to $Q(y)$, which gives the second part of the theorem. ∎

From (12.26), which is the key step in proving that $X^*$ follows the invariant distribution, one can see why it is important to consider sampling that expands backward in time rather than forward. More specifically, consider the coalescence time for the forward chain, letting $\tilde{v}(u_0^\infty)$ be the first index for which

$$\tilde{X}_* := F_{\tilde{v}}(x, u_0^{\tilde{v}})$$

is independent from the starting point, $x$. For any $k \geq 0$, one still has the fact that $F_{\tilde{v}+k}(x, u_0^{\tilde{v}+k})$ does not depend on $x$, but its value depends on $k$ and will not be equal to $\tilde{X}_*$ anymore, which prevents the rest of the proof of theorem 12.9 to carry on.

An equivalent algorithm is described in the next proposition (the proof is easy and left to the reader).

**Proposition 12.10** *Using the same notation as above, the following algorithm generates a perfect sample, $\xi_*$, of the invariant distribution of an ergodic Markov chain.*

*Assume that an infinite sample $u_{-\infty}^0$ of $U$ is available. Given this sequence, the algorithm, starting with $t_0 = 2$, is:*

1. *For all $x \in \mathcal{B}$, define $\xi_{-t}^x, t = -t_0, \ldots, 0$ by $\xi_{-t_0}^x = x$ and $\xi_{-t+1}^x = f(\xi_{-t}^x, u_{-t+1})$.*

2. *If $\xi_0^x$ is constant (independent of $x$), let $\xi_*$ be equal to this constant value and stop. Otherwise, return to step 1 replacing $t_0$ with $2t_0$.*

In practice, the $u_{-k}$'s are only generated when they are needed. But it is important to consider the sequence as fixed: once $u_{-k}$ is generated, it must be stored (or identically regenerated, using the same seed) for further use. It is important to strengthen the fact that this algorithm works backward in time, in the sense that the first states of the sequence are not identical at each iteration, because they are generated using random numbers with indexes further in the past.

Such an algorithm is not feasible when $|\mathcal{B}|$ is too large, since one would have to consider an intractable number of Markov chains (one for each $x \in \mathcal{B}$). However there are cases in which the constancy of $\xi_0^x$ over all $\mathcal{B}$ can be decided from its constancy over a small subset of $\mathcal{B}$.

One situation in which this is true is when the Markov chain is monotone, according to the following definition. Assume that $\mathcal{B}$ can be partially ordered, and that $f$ in (12.24) is increasing in $x$, i.e.,

$$x \leq x' \Rightarrow \forall u, f(x, u) \leq f(x', u). \tag{12.27}$$

Let $\mathcal{B}_{\min}$ and $\mathcal{B}_{\max}$ be the set of minimal and maximal elements in $\mathcal{B}$. Then the sequence coalesces for the algorithm above if and only if it coalesces over $\mathcal{B}_{\min} \cup \mathcal{B}_{\max}$. Indeed, any $x \in \mathcal{B}$ is smaller than some maximal element, and larger than some minimal element in $\mathcal{B}$. By (12.27), these inequalities remain true at each step of the sampling process, which implies that when chains initialized with extremal elements coalesce, so do the other ones. Therefore, it suffices to run the algorithm with extremal configurations only.

One can rewrite (12.27) in terms of transition probabilities $p(x, y)$, assuming that $U$ follows a uniform distribution on $[0, 1]$ and, for all $x \in \mathcal{B}$, there exists a partition $(I_{xy}, y \in \mathcal{B})$ of $\mathcal{B}$, such that
$$f(x, u) = y \Leftrightarrow u \in I_{x,y}$$
and $I_{xy}$ is an interval with length $p_{xy}$. Condition (12.27) is then equivalent to
$$x \leq x' \Rightarrow \forall y \in \mathcal{B}, I_{xy} \subset \bigcup_{y' \geq y} I_{x'y'}.$$
This requires in particular that $\sum_{y \geq y_0} p(x, y) \leq \sum_{y \geq y_0} p(x', y)$ whenever $x \leq x'$ (one says that $p(x, \cdot)$ is stochastically smaller than $p(x', \cdot)$).

One example in which this reduction works is with the ferromagnetic Ising model, for which $\mathcal{B} = \{-1, 1\}^L$ and
$$q(x) = \frac{1}{C} \exp\Big( \sum_{s,t=1, s<t}^{L} \beta_{st} x^{(s)} x^{(t)} \Big)$$
with $\beta_{st} \geq 0$ for all $\{s, t\}$. Then, the Gibbs sampling algorithm iterates the following steps: take a random $s \in \{1, \ldots, L\}$ and update $x^{(s)}$ according to the conditional distribution
$$g_s(y^{(s)} \mid x^{(s^c)}) = \frac{e^{y^{(s)} v_s(x)}}{e^{-v_s(x)} + e^{v_s(x)}}$$
with $v_s(x) = \sum_{t \neq s} \beta_{st} x^{(t)}$. Order $\mathcal{B}$ so that $x \leq \tilde{x}$ if and only if $x^{(s)} \leq \tilde{x}^{(s)}$ for all $s = 1, \ldots, L$. The minimal and maximal elements are unique in this case, with $x_{\min}^{(s)} \equiv -1$ and $x_{\max}^{(s)} \equiv 1$. Moreover, because all $\beta_{st}$ are non-negative, $v_s$ is an increasing function of $x$ so that, if $x \leq \tilde{x}$, then $g_s(1 \mid x^{(s)}) \leq g_s(1 \mid \tilde{x}^{(s)})$.

To define the stochastic iterations, first introduce
$$f_s(x, u) = \begin{cases} 1^{(s)} \wedge x^{(s^c)} & \text{if } u \leq q_s(1 \mid x^{(s)}) \\ (-1)^{(s)} \wedge x^{(s^c)} & \text{if } u > q_s(1 \mid x^{(s)}), \end{cases}$$
which satisfies (12.27). The whole updating scheme can then be implemented with the function
$$f(x, (u, \tilde{u})) = \sum_{s=1}^{L} \delta_{I_s}(\tilde{u}) f_s(x, u)$$

where $(I_s, s \in V)$ is any partition of $[0,1]$ in intervals of length $1/L$. This is still monotonic. The algorithm described in proposition 12.10 can therefore be applied to sample exactly, in finite time, from the ferromagnetic Ising model.

## 12.7 Application: Stochastic approximation with Markovian transitions

Using the material developed in this chapter, we now discuss the convergence of stochastic approximation methods (such as stochastic gradient descent) when the random random variable in the update term follows Markovian transitions. In section 3.3, we considered algorithms in the form

$$\begin{cases} \xi_{t+1} \sim \pi_{X_t} \\ X_{t+1} = X_t + \alpha_{t+1} H(X_t, \xi_{t+1}) \end{cases}$$

where $\xi_t : \Omega \to \mathcal{R}_\xi$ is a random variable. We now want to addres situations in which the random variable $\xi_{t+1}$ is obtained through a transition probability, therefore considering the algorithm

$$\begin{cases} \xi_{t+1} \sim P_{X_t}(\xi_t, \cdot) \\ X_{t+1} = X_t + \alpha_{t+1} H(X_t, \xi_{t+1}) \end{cases} \tag{12.28}$$

Here $P_x$ is, for all $x$, a transition probability from $\mathcal{R}_\xi$ to $\mathcal{R}_\xi$. We will assume that, for all $x \in \mathbb{R}^d$, the Markov chain with transition $P_x$ is geometrically ergodic, and we denote by $\pi_x$ its invariant distribution. We let, as in section 3.3, $\bar{H}(x) = E_{\pi_x}(H(x, \cdot))$. We will use the notation for a function $f : \mathbb{R}^d \times \mathcal{R}_\xi \to \mathbb{R}$

$$P_x f : (x', \xi) \in \mathbb{R}^d \times \mathcal{R}_\xi \mapsto P_x f(x', \xi) = \int_{\mathcal{R}_\xi} f(x', \xi') P_x(\xi, d\xi')$$

and

$$\pi_x f : x' \in \mathbb{R}^d \mapsto \pi_x f(x') = \int_{\mathcal{R}_\xi} f(x', \xi) \pi_x(d\xi).$$

In particular, $\bar{H}(x) = \pi_x H(x)$. We also define $h(x, \xi) = H(x, \xi) - \bar{H}(x)$ and $\tilde{h}(x, \xi) = P_x h(x, \xi)$. We make the following assumptions.

(H1) There exists constants $C_0, C_1, c_2$ such that, for all $x, y \in \mathbb{R}^d$,

$$\sup_{\xi \in \mathcal{R}_\xi} |H(x, \xi)| \leq C_0, \tag{12.29a}$$

$$\sup_{\xi \in \mathcal{R}_\xi} |\tilde{h}(x, \xi)| \leq C_1, \tag{12.29b}$$

$$\sup_{\xi \in \mathcal{R}_\xi} |\tilde{h}(x, \xi) - \tilde{h}(y, \xi)| \leq C_1 |x - y|, \tag{12.29c}$$

$$D_{var}(\pi_x, \pi_y) \leq C_2 |x - y| \tag{12.29d}$$

(H2) There exists $x^* \in \mathbb{R}^d$ and $\mu > 0$ such that, for all $x \in \mathbb{R}^d$

$$(x - x^*)^T \bar{H}(x) \leq -\mu |x - x^*|^2. \tag{12.30}$$

(H3) We assume that there exists a constant $M$ and a non-decreasing function $\rho :$ $[0, +\infty) \rightarrow [0, 1)$ such that, for all probability distributions $Q$ and $Q'$ on $\mathcal{R}_\xi$,

$$D_{\mathrm{var}}(QP_x^n, Q'P_x^n) \leq M\rho(|x|)^n D_{\mathrm{var}}(Q, Q'). \tag{12.31}$$

(H4) The sequence $\alpha_1, \alpha_2, \ldots$ is non-increasing, with

$$\sum_{t=1}^{\infty} \alpha_t = +\infty \quad \text{and} \quad \sum_{t=1}^{\infty} \alpha_t^2 < +\infty. \tag{12.32a}$$

Let $\sigma_t = \sum_{s=1}^{t} \alpha_s$. If $C_1 > 0$, we also require that

$$\lim_{t \to \infty} \alpha_t \sigma_t (1 - \rho(\sigma_t))^{-1} = 0 \tag{12.32b}$$

and

$$\sum_{s=2}^{t} \alpha_s^2 \sigma_s (1 - \rho(\sigma_s))^{-2} < \infty. \tag{12.32c}$$

Given this, the following theorem holds.

**Theorem 12.11** *Assuming (H1) to (H4), the sequence defined by (12.28) is such that*

$$\lim_{t \to \infty} \mathbb{E}(|X_t - x_*|^2) = 0$$

**Remark 12.12** Condition (H1) assumes that $H$ is bounded and uniformly Lipschitz in $x$, which is more restrictive than what was assumed in section 3.3.2, but applies, for example, to situations considered in Younes [206] and later in this book in section 17.2.2.

Condition (H3) implies that the Markov chain with transition $P_x$ is uniformly geometrically ergodic, but the ergodicity rate may depend on $x$ and in particular converge to 1 when $x$ tends to $\infty$, which is the situation targeted in this theorem.

The reader may refer to [208] for a general discussion of this problem with relaxed hypotheses and almost sure convergence, at the expense of significantly longer proofs. ◆

PROOF We note that, from (12.29a), one has

$$|X_t - x_*| \leq C_0 \sigma_t |X_0 - x_*|. \tag{12.33}$$

Similarly to section 3.3.2, we let $A_t = |X_t - x_*|^2$ and $a_t = \mathbb{E}(A_t)$. One can then write

$$A_{t+1} = A_t + 2\alpha_{t+1}(X_t - x_*)^T \bar{H}(X_t) + 2\alpha_{t+1}(X_t - x_*)^T(H(X_t, \xi_{t+1}) - \bar{H}(X_t)) + \alpha_{t+1}^2 |H(X_t, \xi_{t+1})|^2$$

but we do not have

$$\mathbb{E}((X_t - x_*)^T(H(X_t, \xi_{t+1}) - \bar{H}(X_t)) \,|\, \mathcal{U}_t) = 0$$

anymore, where $\mathcal{U}_t$ is the $\sigma$-algebra of all past events up to time $t$ (all events depending of $X_s, \xi_s$, $s \leq t$). Indeed the Markovian assumption implies that

$$\mathbb{E}((X_t - x_*)^T(H(X_t, \xi_{t+1}) - \bar{H}(X_t)) \,|\, \mathcal{U}_t) = (X_t - x_*)^T \left( \int_{\mathcal{R}_\xi} H(X_t, \xi) P_{X_t}(\xi_t, d\xi) - \bar{H}(X_t) \right)$$
$$= (X_t - x_*)^T ((P_{X_t} H(X_t, \cdot))(\xi_t) - \bar{H}(X_t)),$$

which does not vanish in general. Following Benveniste et al. [25], this can be addressed by introducing the solution $g(x, \cdot)$ of the "Poisson equation"

$$g(x, \cdot) - P_x g(x, \cdot) = h(x, \cdot). \tag{12.34}$$

(Recall that $h(x, \xi) = H(x, \xi) - \bar{H}(x)$.) One can then write

$$(X_t - x_*)^T h(X_t, \xi_{t+1}) = (X_t - x_*)^T (g(X_t, \xi_{t+1}) - P_{X_t} g(X_t, \xi_{t+1})$$

and

$$A_{t+1} \leq (1 - 2\alpha_{t+1}\mu)A_t + 2\alpha_{t+1}(X_t - x_*)^T(g(X_t, \xi_{t+1}) - P_{X_t} g(X_t, \xi_t)))$$
$$+ 2\alpha_{t+1}(X_t - x_*)^T P_{X_t} g(X_t, \xi_t) - 2\alpha_{t+1}(X_t - x_*)^T P_{X_t} g(X_t, \xi_{t+1}) + \alpha_{t+1}^2 |H(X_t, \xi_{t+1})|^2$$

Introducing the notation

$$\eta_{st} = \mathbb{E}((X_s - x_*)^T P_{X_s} g(X_s, \xi_t)).$$

Using the fact that

$$\mathbb{E}\left((X_t - x_*)^T(g(X_t, \xi_{t+1}) - P_{X_t} g(X_t, \xi_t))) \,|\, \mathcal{U}_t\right) = 0$$

and and noting that $|H(X_t, \xi_{t+1})|^2 \leq C_0^2$, this gives, after taking expectations,

$$a_{t+1} \leq (1 - 2\alpha_{t+1}\mu)a_t + 2\alpha_{t+1}\eta_{tt} - 2\alpha_{t+1}\eta_{t,t+1} + \alpha_{t+1}^2 C_0^2.$$

Applying lemma 3.25, and letting $v_{s,t} = \prod_{j=s+1}^t (1 - 2\alpha_{j+1}\mu)$, we get

$$a_t \leq a_0 v_{0,t} + 2 \sum_{s=1}^t v_{s,t}\alpha_{s+1}(\eta_{ss} - \eta_{s,s+1}) + C_0^2 \sum_{s=1}^t v_{s,t}\alpha_{s+1}^2.$$

We now want to ensure that each term in the upper bound converges to 0. Similarly to section 3.3.2, (12.32a) implies that this holds the first and last terms and we therefore focus on the middle one, writing

$$\sum_{s=1}^{t} v_{s,t}\alpha_{s+1}(\eta_{ss} - \eta_{s,s+1}) = v_{1,t}\alpha_2\eta_{11} - \alpha_{t+1}\eta_{t,t+1} + \sum_{s=2}^{t}(v_{s,t}\alpha_{s+1} - v_{s-1,t}\alpha_s)\eta_{ss} \quad (12.35)$$

$$+ \sum_{s=2}^{t} v_{s-1,t}\alpha_s(\eta_{ss} - \eta_{s-1,s})$$

We will need the following estimates on the function $g$ in (12.34), which is defined by

$$g(x,\xi) = \sum_{n=0}^{\infty} P_x^n h(x,\xi) = h(x,\xi) + \sum_{n=0}^{\infty} P_x^n \tilde{h}(x,\xi).$$

**Lemma 12.13** *We have*

$$|g(x,\cdot)| \le C_0 + 2C_1 M(1 - \rho(x))^{-1}, \quad (12.36a)$$

$$|P_x g(x,\cdot)| \le 2C_1 M(1 - \rho(x))^{-1}. \quad (12.36b)$$

*and, for all $x, y \in \mathbb{R}^d$ and $\xi \in \mathcal{R}_\xi$*

$$|P_x g(x,\xi) - P_y(g(y,\xi)| = M^2 C_1 C_2(1 - \bar{\rho})^{-2} + MC_1(1 + C_2)(1 - \bar{\rho})^{-1}. \quad (12.37)$$

*with $\bar{\rho} = \max(\rho(|x|), \rho(|y|))$.*

Using lemma lemma 12.13 (which is proved at the end of the section), we can control the terms intervening in (12.35). Note that the first term, $v_{1t}\alpha_2\eta_{11}$, converges to 0 since (12.32a) implies that $v_{1t}$ converges to 0.

We have,

$$\alpha_{t+1}|\mathbb{E}((X_t - x_*)^T P_{X_t} g(X_t, \xi_{t+1}))| \le 2MC_1\alpha_{t+1}\sigma_t(1 - \rho(\sigma_t))^{-1},$$

so that (12.32b) implies that $\alpha_{t+1}\eta_{t,t+1} \to 0$.

and since $\alpha_{s+1} \le \alpha_s$, we have

$$\left|\sum_{s=2}^{t}(v_{s,t}\alpha_{s+1} - v_{s-1,t}\alpha_s)\eta_{ss}\right| \le \sum_{s=2}^{t}|v_{st}\alpha_s - v_{s,t}\alpha_{s+1}||\eta_{ss}|$$

$$\le MC_1 \sum_{s=2}^{t}|v_{s-1,t}\alpha_s - v_{s,t}\alpha_{s+1}|\alpha_{s+1}\sigma_s(1 - \rho(\sigma_s))^{-1}$$

$$\le C \sum_{s=2}^{t}|v_{s-1,t}\alpha_s - v_{s,t}\alpha_{s+1}|$$

for some constant $C$, since $\alpha_{s+1}\sigma_s(1-\rho(\sigma_s))^{-1}$ is bounded. Writing

$$v_{s,t}\alpha_{s+1} - v_{s-1,t}\alpha_s = v_{st}(\alpha_{s+1} - \alpha_s + 2\mu\alpha_s^2),$$

we get (using $\alpha_{s+1} \leq \alpha_s$)

$$\sum_{s=2}^{t} |v_{s-1,t}\alpha_s - v_{s,t}\alpha_{s+1}| \leq \sum_{s=2}^{t} v_{st}(\alpha_s - \alpha_{s+1}) + \sum_{s=2}^{t} v_{st}2\mu\alpha_s^2.$$

Since both $\sum_s(\alpha_s - \alpha_{s+1})$ and $\sum_{s=2}^{t}\alpha_s^2$ converge (the former is just $\alpha_1$), lemma 3.26 implies that

$$\sum_{s=2}^{t}(v_{s,t}\alpha_{s+1} - v_{s-1,t}\alpha_s)\eta_{ss}$$

tends to zero. The last term to consider is

$$\sum_{s=2}^{t} v_{s-1,t}\alpha_s(\eta_{ss} - \eta_{s-1,s}) = \sum_{s=2}^{t} v_{s-1,t}\alpha_s\mathbb{E}((X_s - X_{s-1})^T P_{X_s}g(X_s, \xi_s))$$

$$+ \sum_{s=2}^{t} v_{s-1,t}\alpha_s\mathbb{E}((X_{s-1} - x_*)^T(P_{X_s}g(X_s, \xi_s)) - P_{X_{s-1}}g(X_{s-1}, \xi_s))).$$

We have

$$\left| \sum_{s=2}^{t} v_{s-1,t}\alpha_s\mathbb{E}((X_s - X_{s-1})^T P_{X_s}g(X_s, \xi_s)) \right| \leq 2C_0C_1M\sum_{s=2}^{t} v_{s-1,t}\alpha_s^2(1-\rho(\sigma_s))^{-1}$$

and

$$\left| \sum_{s=2}^{t} v_{s-1,t}\alpha_s\mathbb{E}((X_{s-1} - x_*)^T(P_{X_s}g(X_s, \xi_s)) - P_{X_{s-1}}g(X_{s-1}, \xi_s))) \right|$$

$$\leq 2M^2C_0C_1(1+C_2)|X_0 - x_*|\sum_{s=2}^{t} v_{s-1,t}\alpha_s^2\sigma_s(1-\rho(\sigma_s))^{-2}$$

and lemma 3.26 implies that both terms vanish at infinity. This concludes the proof of theorem 12.11. ∎

PROOF (PROOF OF LEMMA 12.13) Condition (H3) and proposition 12.3 and imply that (since $\pi_x\tilde{h} = 0$)

$$|P_x^n\tilde{h}(x, \xi)| \leq D_{var}(P_x^n(\xi, \cdot), \pi_x)osc(\tilde{h}(x, \cdot)) \leq 2C_1M\rho(x)^n$$

so that $g$ is well defined with

$$|g(x, \cdot)| \leq C_0 + 2C_1M(1-\rho(x))^{-1},$$
$$|P_xg(x, \cdot)| \leq 2C_1M(1-\rho(x))^{-1}.$$

We will also need to control differences of the kind

$$P_x g(x, \xi) - P_y g(y, \xi).$$

We consider the $n$th term in the series, writing

$$P_x^n \tilde{h}(x, \xi) - P_y^n \tilde{h}(y, \xi) = \sum_{k=0}^{n-1} (P_x^{n-k} P_y^k \tilde{h}(y, \xi) - (P_x^{n-k-1} P_y^{k+1} \tilde{h}(y, \xi))$$
$$+ P_x^n \tilde{h}(x, \xi) - P_x^n \tilde{h}(y, \xi).$$

This gives

$$P_x^n \tilde{h}(x, \xi) - P_y^n \tilde{h}(y, \xi) = \sum_{k=0}^{n-1} P_x^{n-k-1} (P_x P_y^k \tilde{h}(y, \xi) - P_y^{k+1} \tilde{h}(y, \xi) - \pi_x P_y^k \tilde{h}(y) + \pi_x P_y^{k+1} \tilde{h}(y))$$

$$+ \sum_{k=0}^{n-1} (\pi_x P_y^k \tilde{h}(y) - \pi_x P_y^{k+1} \tilde{h}(y)) + P_x^n \tilde{h}(x, \xi) - P_x^n \tilde{h}(y, \xi)$$

$$= \sum_{k=0}^{n-1} P_x^{n-k-1} (P_x P_y^k \tilde{h}(y, \xi) - P_y^{k+1} \tilde{h}(y, \xi) - \pi_x P_y^k \tilde{h}(y) + \pi_x P_y^{k+1} \tilde{h}(y))$$

$$+ \pi_x \tilde{h}(y) - \pi_x P_y^n \tilde{h}(y) + P_x^n \tilde{h}(x, \xi) - P_x^n \tilde{h}(y, \xi)$$

Finally

$$P_x^n h(x, \xi) - P_y^n h(y, \xi) = \sum_{k=0}^{n-1} P_x^{n-k-1} (P_x P_y^k \tilde{h}(y, \xi) - P_y^{k+1} \tilde{h}(y, \xi) - \pi_x P_y^k \tilde{h}(y) + \pi_x P_y^{k+1} \tilde{h}(y))$$

$$+ P_x^n (\tilde{h}(x, \xi) - \tilde{h}(y, \xi) + \pi_x \tilde{h}(y)) - (\pi_x - \pi_y) P_y^n \tilde{h}(y)$$

Using proposition 12.3, we can write, letting $\bar{\rho} = \max(\rho(|x|), \rho(|y|))$,

$$|P_x^{n-k-1} (P_x P_y^k \tilde{h}(y, \xi) - P_y^{k+1} \tilde{h}(y, \xi) - \pi_x P_y^k \tilde{h}(y, \xi) + \pi_x P_y^{k+1} \tilde{h}(y, \xi))|$$
$$\leq M \bar{\rho}^{n-k-1} osc(P_x P_y^k \tilde{h}(y, \xi) - P_y^{k+1} \tilde{h}(y, \xi))$$
$$\leq C_2 M \bar{\rho}^{n-k-1} |x - y| osc(P_y^k \tilde{h}(y, \xi))$$
$$\leq C_2 C_1 M^2 \bar{\rho}^{n-1} |x - y|$$

We also have
$$|P_x^n (\tilde{h}(x, \xi) - \tilde{h}(y, \xi) + \pi_x \tilde{h}(y, \xi))| \leq M C_1 \bar{\rho}^n |x - y|$$
and
$$|(\pi_x - \pi_y) P_y^n \tilde{h}(y, \xi)| \leq M C_2 C_1 \bar{\rho}^n |x - y|$$

so that

$$|P_x^n h(x,\xi) - P_y^n h(y,\xi)| \leq M C_1 \bar{\rho}^{n-1}(n M C_2 + (1 + C_2)\bar{\rho})|x - y|$$

From this, it follows that

$$|P_x g(x,\xi) - P_y(g(y,\xi)| \leq M C_1 \sum_{n=1}^{\infty} \bar{\rho}^{n-1}(n M C_2 + (1 + C_2)|x - y|$$

$$= M^2 C_1 C_2 (1 - \bar{\rho})^{-2} + M C_1 (1 + C_2)(1 - \bar{\rho})^{-1}. \qquad \blacksquare$$

# Chapter 13

# Markov Random Fields

With this chapter, we start a discussion of large-scale statistical models in data science, starting with graphical models (Markov random fields and Bayesian networks) before discussing more recent approaches using, notably, deep learning. Important textbook references for the present chapter include Pearl [151], Ancona et al. [8], Winkler [203], Lauritzen [114], Cowell et al. [56], Koller and Friedman [108].

## 13.1 Independence and conditional independence

### 13.1.1 Definitions

We consider random variables $X, Y, Z \ldots$, and denote by $\mathcal{R}_X, \mathcal{R}_Y, \mathcal{R}_Z \ldots$ the sets in which they take their values. We discuss in this section concepts of independence and conditional independence between random variables. To simplify the exposition, we will work (unless mentioned otherwise) with discrete random variables ($X$ is discrete if $\mathcal{R}_X$ is finite or countable)[1]. We start with a basic definition.

**Definition 13.1** *Two discrete random variables $X : \Omega \to \mathcal{R}_X$ and $Y : \Omega \to \mathcal{R}_Y$ are independent if and only if*

$$\forall x \in \mathcal{R}_X, \forall y \in \mathcal{R}_Y : \mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x)\mathbb{P}(Y = y).$$

The general definition for arbitrary r.v.'s is that

$$\mathbb{E}(f(X)g(Y)) = \mathbb{E}(f(X))\mathbb{E}(g(Y))$$

for any pair of (measurable) non-negative functions $f : \mathcal{R}_X \to [0, +\infty)$ and $g : \mathcal{R}_Y \to [0, +\infty)$.

---

[1]In the general case, $\mathcal{R}_X, \mathcal{R}_Y, \ldots$ are metric spaces with a countable dense subset with $\sigma$-algebras $\mathcal{S}_X, \mathcal{S}_Y, \ldots$

One can easily check that $X$ and $Y$ are independent if and only if, for any non-negative function $g : \mathcal{R}_Y \to \mathbb{R}$, one has

$$\mathbb{E}(g(Y) \mid X) = \mathbb{E}(g(Y)).$$

**Notation 13.2** Independence is a property that involves two variables $X$ and $Y$ and an underlying probability distribution $\mathbb{P}$. Independence of $X$ and $Y$ relative to $\mathbb{P}$ will be denoted $(X \perp\!\!\!\perp Y)_{\mathbb{P}}$. However we will only write $X \perp\!\!\!\perp Y$ when there is no ambiguity on $\mathbb{P}$. ♦

More than independence, the concept of conditional independence will be fundamental in this chapter. It requires three variables, say $X, Y, Z$. Returning to the discrete case, one says that $X$ and $Y$ are conditionally independent given $Z$ is, for any $x \in \mathcal{R}_X$, $y \in \mathcal{R}_Y$ and $z \in \mathcal{R}_Z$ such that $\mathbb{P}(Z = z) > 0$,

$$\mathbb{P}(X = x, Y = y \mid Z = z) = \mathbb{P}(X = x \mid Z = z)\,\mathbb{P}(Y = y \mid Z = z). \tag{13.1}$$

An equivalent statement is that, for any $z$ such that $\mathbb{P}(Z = z) \neq 0$, $X$ and $Y$ are independent when $\mathbb{P}$ is replaced by the conditional distribution $\mathbb{P}(\cdot \mid Z = z)$.

In the general case conditional independence means that, for any pair of non-negative measurable functions $f$ and $g$,

$$\mathbb{E}(f(X)g(Y) \mid Z) = \mathbb{E}(f(X) \mid Z)\,\mathbb{E}(g(Y) \mid Z). \tag{13.2}$$

From now, we restrict our discussion to discrete random variables.

Multiplying both terms in (13.1) by $\mathbb{P}(Z = z)^2$, we get the equivalent statement: $X$ and $Y$ are conditionally independent given $Z$ if and only if,

$$\forall x, y, z : \mathbb{P}(X = x, Y = y, Z = z)\mathbb{P}(Z = z) = \mathbb{P}(X = x, Z = z)\,\mathbb{P}(Y = y, Z = z). \tag{13.3}$$

Note that the identity is meaningful, and always true, for $\mathbb{P}(Z = z) = 0$, so that this case does not need to be excluded anymore.

Conditional independence can be interpreted by the statement that $X$ brings no more information on $Y$ than what is already provided by $Z$: one has

$$\mathbb{P}(Y = y \mid X = x, Z = z) = \frac{\mathbb{P}(Y = y, X = x, Z = z)}{\mathbb{P}(X = x, Z = z)} = \frac{\mathbb{P}(Y = y, Z = z)}{\mathbb{P}(Z = z)}$$

as directly deduced from (13.3). (This computation being valid as soon as $\mathbb{P}(X = x, Z = z) > 0$.)

**Notation 13.3** To indicate that $X$ and $Y$ are conditionally independent given $Z$ for the distribution $\mathbb{P}$, we will write $(X \perp\!\!\!\perp Y \mid Z)_{\mathbb{P}}$ or simply $(X \perp\!\!\!\perp Y \mid Z)$. ♦

So we have the equivalence:

$$(X \perp\!\!\!\perp Y \mid Z)_{\mathbb{P}} \Leftrightarrow \left( \forall z : \mathbb{P}(Z = z) > 0 \Rightarrow (X \perp\!\!\!\perp Y)_{\mathbb{P}(\cdot \mid Z=z)} \right).$$

Absolute independence is like "independence conditional to no variable", and we will use the notation $\emptyset$ for the "empty" random variable that contains no information (for example, a set-valued random variable that always returns the empty set, or any constant variable). So we have the tautology

$$X \perp\!\!\!\perp Y \Leftrightarrow (X \perp\!\!\!\perp Y \mid \emptyset).$$

Note that, dealing with discrete variables, all previous definitions automatically extend to groups of variables: for example, if $Z_1$, $Z_2$ are two discrete variables, so is $Z = (Z_1, Z_2)$ and we immediately obtain a definition for the conditional independence of $X$ and $Y$ given $Z_1$ and $Z_2$, denoted $(X \perp\!\!\!\perp Y \mid Z_1, Z_2)$.

### 13.1.2 Fundamental properties

Proposition 13.5 below lists important properties of conditional independence that will be used repeatedly in this chapter. Before stating this proposition, we need the following definition.

**Definition 13.4** *One says that the joint distribution of the random variables $(X_1, \ldots, X_N)$ is positive if there exists subsets $\tilde{R}_k \subset \mathcal{R}_{X_k}$, $k = 1, \ldots, N$ such that $\mathbb{P}(X_k \in \tilde{R}_k) = 1$ and:*

$$\mathbb{P}(X_1 = x_1, \ldots, X_N = x_N) > 0$$

*if $x_k \in \tilde{R}_k$, $k = 1, \ldots, N$.*

Note that the condition implies $\mathbb{P}(X_k = x_k) > 0$ for all $x_k \in \tilde{R}_k$, so that $\tilde{R}_k = \{x_k \in \mathcal{R}_{X_k} : \mathbb{P}(X_k = x_k) > 0\}$, i.e., $\tilde{R}_k$ is the support of $P_{X_k}$. One can interpret the definition as expressing the fact that any conjunction of events for different $X_k$'s has positive probability, as soon as each of them has positive probability (if all events may occur, then they may occur together).

Note that the sets $\tilde{R}_k$ depend on $X_1, \ldots, X_N$. However, if this family of variables is fixed, there is no loss in generality in restricting the space $\mathcal{R}_{X_k}$ to $\tilde{R}_k$ and there for assume that $\mathbb{P}(X_1 = x_1, \ldots, X_N = x_N) > 0$ everywhere.

**Proposition 13.5** *Let $X, Y, Z$ and $W$ be random variables. The following properties are true.*

*(CI1) Symmetry: $(X \perp\!\!\!\perp Y \mid Z) \Rightarrow (Y \perp\!\!\!\perp X \mid Z)$.*

*(CI2) Decomposition:* $(X \perp\!\!\!\perp (Y, W) \mid Z) \Rightarrow (X \perp\!\!\!\perp Y \mid Z)$.

*(CI3) Weak union:* $(X \perp\!\!\!\perp (Y, W) \mid Z) \Rightarrow (X \perp\!\!\!\perp Y \mid (Z, W))$.

*(CI4) Contraction:* $(X \perp\!\!\!\perp Y \mid Z)$ *and* $(X \perp\!\!\!\perp W \mid (Z, Y)) \Rightarrow (X \perp\!\!\!\perp (Y, W) \mid Z)$.

*(CI5) Intersection: assume that the joint distribution of* $W, Y$ *and* $Z$ *is positive. Then*

$$(X \perp\!\!\!\perp W \mid (Z, Y)) \text{ and } (X \perp\!\!\!\perp Y \mid (Z, W)) \Rightarrow (X \perp\!\!\!\perp (Y, W) \mid Z).$$

PROOF Properties (CI1) and (CI2) are easily deduced from (13.3) and left to the reader. To prove the last three, we will use the notation $P(x), P(x, y)$ etc. instead of $\mathbb{P}(X = x), \mathbb{P}(X = x, Y = y)$, etc. to save space. Identities are assumed to hold for all $x, y, z, w$ unless stated otherwise.

For (CI3), we must prove, according to (13.3), that

$$P(x, y, z, w)P(z, w) = P(x, z, w)P(y, z, w) \tag{13.4}$$

whenever $P(x, y, z, w)P(z) = P(x, z)P(y, z, w)$. Summing this last equation over $y$ (or applying (CI2)) yields $P(x, z, w)P(z) = P(x, z)P(z, w)$. We can note that all terms in (13.4) vanish when $P(z) = 0$, so that the identity is true in this case. When $P(z) \neq 0$, the right-hand side of (13.4) becomes

$$(P(x, z)P(z, w)/P(z))P(y, z, w) = (P(x, z)P(y, z, w)/P(z))P(z, w) = P(x, y, z, w)P(z, w),$$

using once again the hypothesis. This proves (CI3).

For (CI4), the hypotheses are

$$\begin{cases} P(x, y, z)P(z) = P(x, z)P(y, z) \\ P(x, y, z, w)P(y, z) = P(x, y, z)P(y, z, w) \end{cases}$$

and the conclusion must be

$$P(x, y, z, w)P(z) = P(x, z)P(y, z, w). \tag{13.5}$$

Since (13.5) is true when $P(y, z) = 0$, we assume that this probability does not vanish and write

$$\begin{aligned} P(x, y, z, w)P(z) &= P(x, y, z)P(z)P(y, z, w)/P(y, z) \\ &= P(x, z)P(y, z)P(y, z, w)/P(y, z) \\ &= P(x, z)P(y, z, w) \end{aligned}$$

yielding (13.5).

For (CI5), assuming

$$\begin{cases} P(x,y,z,w)P(y,z) = P(x,y,z)P(y,z,w) \\ P(x,y,z,w)P(z,w) = P(x,z,w)P(y,z,w), \end{cases} \tag{13.6}$$

we want to show that

$$P(x,y,z,w)P(z) = P(x,z)P(y,z,w).$$

Since this identity is true when any of the events $W = w, Y = y$ or $Z = z$ has zero probability, we can assume that their probabilities are positive, which, by assumption, also implies that all joint probabilities are positive. From the two identities, we get

$$P(x,y,z,w)/P(y,z,w) = P(x,y,z)/P(y,z) = P(x,z,w)/P(z,w)$$

This implies

$$P(x,y,z) = P(y,z)P(x,z,w)/P(z,w)$$

that we can sum over $y$ to obtain

$$P(x,z) = P(z)P(x,z,w)/P(z,w)$$

We therefore get

$$P(x,y,z,w)/P(y,z,w) = P(x,z,w)/P(z,w) = P(x,z)/P(z),$$

which is what we wanted. ∎

A counter-example of (CI5) when the positivity assumption is not satisfied can be built as follows: let $X$ be a Bernoulli random variable, and let $Y = W = X$. Let $Z$ be any Bernoulli random variable, independent from $X$. Given $Z$ and $W$, $X$ and $Y$ are constant and therefore independent. Similarly, given $Z$ and $Y$, $X$ and $W$ are constant and therefore independent. However, given $Z$, $X$ and $(Y, W)$ are not independent (they are equal and non constant).

### 13.1.3  Mutual independence

Another concept of interest is the mutual (conditional) independence of more than two random variables. The random variables $(X_1, \ldots, X_n)$ are mutually conditionally independent given $Z$ if and only if

$$\mathbb{E}(f_1(X_1) \cdots f_n(X_n) \mid Z) = \mathbb{E}(f_1(X_1) \mid Z) \cdots \mathbb{E}(f_n(X_n) \mid Z)$$

for any non-negative measurable functions $f_1, \ldots, f_n$. In terms of discrete probabilities, this can be written as

$$P(X_1 = x_1, \ldots, X_n = x_n, Z = z)P(Z = z)^{n-1} =$$
$$P(X_1 = x_1, Z = z) \cdots P(X_n = x_n, Z = z).$$

This will be summarized with the notation

$$(X_1 \perp\!\!\!\perp \cdots \perp\!\!\!\perp X_n \mid Z).$$

We have the proposition

**Proposition 13.6** *For variables $X_1, \ldots, X_n$ and $Z$, the following properties are equivalent.*

*(i)  $(X_1 \perp\!\!\!\perp \cdots \perp\!\!\!\perp X_n \mid Z)$;*

*(ii)  For all $S, T \subset \{1, \ldots, n\}$ with $S \cap T = \emptyset$, we have: $((X_i, i \in S) \perp\!\!\!\perp (X_j, j \in T) \mid Z)$;*

*(iii)  For all $s \in \{1, \ldots, n\}$, we have: $(X_s \perp\!\!\!\perp (X_t, t \neq s) \mid Z)$;*

*(iv)  For all $s \in \{2, \ldots, n\}$, we have: $(X_s \perp\!\!\!\perp (X_1, \ldots, X_{s-1}) \mid Z)$.*

PROOF  It is clear that (i) $\Rightarrow \cdots \Rightarrow$ (iv) so it suffices to prove that (iv) $\Rightarrow$ (i). For this, simply write (applying (iv) repeatedly to $s = n-1, n-2, \ldots$)

$$
\begin{aligned}
\mathbb{E}(f_1(X_1) \cdots f_n(X_n) \mid Z) &= \mathbb{E}(f_1(X_1) \cdots f_{n-1}(X_{n-1}) \mid Z) \mathbb{E}(f_n(X_n) \mid Z) \\
&= \mathbb{E}(f_1(X_1) \cdots f_{n-2}(X_{n-2}) \mid Z) \mathbb{E}(f_{n-1}(X_{n-1}) \mid Z) \\
&\qquad \mathbb{E}(f_n(X_n) \mid Z) \\
&\vdots \\
&= \mathbb{E}(f_1(X_1) \mid Z) \cdots \mathbb{E}(f_n(X_n) \mid Z).
\end{aligned}
$$

### 13.1.4   Relation with Information Theory

Several concepts in information theory are directly related to independence between random variables. Recall that the (Shannon) entropy of a discrete probability distribution over a finite set $\mathcal{R}$ is defined by

$$\mathcal{H}(P) = -\sum_{\omega \in \mathcal{R}} \log P(\omega) P(\omega). \tag{13.7}$$

Similarly, the entropy of a random variable $X : \Omega \to \mathcal{R}_X$ is defined by

$$\mathcal{H}(X) \overset{\Delta}{=} \mathcal{H}(P_X) = -\sum_{x \in \mathcal{R}_X} \log P(X = x) P(X = x). \tag{13.8}$$

The entropy is always non-negative, and provides a measure of the uncertainty associated to $P$. For a given finite set $\mathcal{R}$, it is maximal when $P$ is uniform over $\mathcal{R}$, and minimal (and vanishes) when $P$ is supported by a single $\omega \in \mathcal{R}$ (i.e. $P(\omega) = 1$).

One defines the entropy of two or more random variables as the entropy of their joint distribution, so that, for example,

$$\mathcal{H}(X, Y) = - \sum_{(x,y)\in\mathcal{R}_X\times\mathcal{R}_Y} \log \mathbb{P}(X = x, Y = y)\mathbb{P}(X = x, Y = y).$$

We have the proposition:

**Proposition 13.7** *For random variables $X_1,\dots,X_n$, one has*

$$\mathcal{H}(X_1,\dots,X_n) \le \mathcal{H}(X_1) + \cdots + \mathcal{H}(X_n)$$

*with equality if and only if $(X_1,\dots,X_n)$ are mutually independent.*

Proof The proof of this proposition uses properties of the Kullback-Leibler divergence (c.f. (4.3)), given by, for two probability distributions $\pi$ and $\pi'$ on a finite set $\mathcal{B}$,

$$KL(\pi\|\pi') = \sum_{\omega\in\mathcal{B}} \pi(\omega)\log\frac{\pi(\omega)}{\pi'(\omega)}.$$

with the convention $\pi\log(\pi/\pi') = 0$ if $\pi = 0$ and $= \infty$ if $\pi > 0$ and $\pi' = 0$. Returning to proposition 13.7, a straightforward computation (which is left to the reader) shows that

$$\mathcal{H}(X_1) + \cdots + \mathcal{H}(X_n) - \mathcal{H}(X_1,\dots,X_n) = KL(\pi\|\pi')$$

with $\pi(x_1,\dots,x_n) = \mathbb{P}(X_1 = x_1,\dots,X_n = x_n)$ and $\pi'(x_1,\dots,x_n) = \prod_{k=1}^n \mathbb{P}(X_k = x_k)$. This makes proposition 13.7 a direct consequence of proposition 4.1. ∎

The mutual information between two random variables $X$ and $Y$ is defined by

$$\mathcal{I}(X, Y) = \mathcal{H}(X) + \mathcal{H}(Y) - \mathcal{H}(X, Y). \tag{13.9}$$

From proposition 13.7, $\mathcal{I}(X, Y)$ is nonnegative and vanishes if and only if $X$ and $Y$ are independent. Also from the proof of proposition 13.7, $\mathcal{I}(X, Y)$ is equal to $KL(P_{(X,Y)}\|P_X\otimes P_Y)$ where the first probability is the joint distribution of $X$ and $Y$ and the second one the product of the marginals of $X$ and $Y$, which coincides with $P_{X,Y}$ if and only if $X$ and $Y$ are independent.

If $X$ and $Y$ are two random variables, and $y \in \mathcal{R}_Y$ with $\mathbb{P}(Y = y) > 0$, the entropy of the conditional probability $x \mapsto \mathbb{P}(X = x \mid Y = y)$ is denoted $\mathcal{H}(X \mid Y = y)$, and is a function of $y$. The conditional entropy of $X$ given $Y$, denoted $\mathcal{H}(X \mid Y)$ is the expectation of $\mathcal{H}(X \mid Y = y)$ for the distribution of $Y$, i.e.,

$$\mathcal{H}(X \mid Y) = \sum_{y\in\mathcal{R}_Y} \mathcal{H}(X \mid Y = y)\mathbb{P}(Y = y)$$

$$= - \sum_{x\in R_X}\sum_{y\in\mathcal{R}_Y} \log \mathbb{P}(X = x \mid Y = y)\mathbb{P}(X = x, Y = y).$$

So, we have (with a straightforward proof)

**Proposition 13.8** *Given two random variables $X$ and $Y$, we have*

$$
\begin{aligned}
\mathcal{H}(X \mid Y) &= -E_{X,Y}(\log \mathbb{P}(X = \cdot \mid Y = \cdot)) && \text{(13.10)}\\
&= \mathcal{H}(X, Y) - \mathcal{H}(Y)
\end{aligned}
$$

This proposition also immediately yields:

$$
\mathcal{I}(X, Y) = \mathcal{H}(X) - \mathcal{H}(X \mid Y) = \mathcal{H}(Y) - \mathcal{H}(Y \mid X). \tag{13.11}
$$

The identity $\mathcal{H}(X, Y) = \mathcal{H}(X \mid Y) + \mathcal{H}(Y)$ that is deduced from proposition 13.8 can be generalized to more than two random variables (the proof being left to the reader), yielding, if $X_1, \ldots, X_n$ are random variables:

$$
\mathcal{H}(X_1, \ldots, X_n) = \sum_{k=1}^{n} \mathcal{H}(X_k \mid X_1, \ldots, X_{k-1}). \tag{13.12}
$$

If $Z$ is an additional random variable, the following identity is obtained by applying the previous one to conditional distributions given $Z = z$ and taking averages over $z$:

$$
\mathcal{H}(X_1, \ldots, X_n \mid Z) = \sum_{k=1}^{n} \mathcal{H}(X_k \mid X_1, \ldots, X_{k-1}, Z). \tag{13.13}
$$

The following proposition characterizes conditional independence in terms of entropy.

**Proposition 13.9** *Let $X, Y$ and $Z$ be three random variables. The following statements are equivalent.*

*(i)  $X$ and $Y$ are conditionally independent given $Z$.*

*(ii)  $\mathcal{H}(X, Y \mid Z) = \mathcal{H}(X \mid Z) + \mathcal{H}(Y \mid Z)$*

*(iii)  $\mathcal{H}(X \mid Y, Z) = \mathcal{H}(X \mid Y)$*

*Moreover, when (i) to (iii) are satisfied, we have:*

*(iv)  $\mathcal{I}(X, Y) \leq \min(\mathcal{I}(X, Z), \mathcal{I}(Y, Z))$.*

Proof From proposition 13.7, we have, for any three random variables $X, Y, Z$, and any $z$ such that $P(Z = z) > 0$,

$$
\mathcal{H}(X, Y \mid Z = z) \leq \mathcal{H}(X \mid Z = z) + \mathcal{H}(Y \mid Z = z).
$$

Taking expectations on both sides implies the important inequality

$$\mathcal{H}(X, Y \mid Z) \leq \mathcal{H}(X \mid Z) + \mathcal{H}(Y \mid Z) \tag{13.14}$$

and equality occurs if and only if $\mathbb{P}(X = x, Y = y \mid Z = z) = \mathbb{P}(X = x \mid Z = z)\mathbb{P}(Y = y \mid Z = z)$ whenever $\mathbb{P}(Z = z) > 0$, that is, if and only if $X$ and $Y$ are conditionally independent given $Z$. This proves that (i) and (ii) are equivalent. The fact that (ii) and (iii) are equivalent comes from (13.13), which gives, for any three random variables

$$\mathcal{H}(X, Y \mid Z) = \mathcal{H}(X \mid Y, Z) + \mathcal{H}(Y \mid Z). \tag{13.15}$$

To prove that (i)-(iii) implies (iv), we note that (13.14) and (13.15) imply that, for any three random variables:

$$\mathcal{H}(X \mid Y, Z) \leq \mathcal{H}(X \mid Y).$$

If $X$ and $Y$ are conditionally independent given $Z$, then the right-hand side is equal to $\mathcal{H}(X \mid Z)$ and this yields

$$\mathcal{I}(X, Y) = \mathcal{H}(X) - \mathcal{H}(X \mid Y) \leq \mathcal{H}(X) - \mathcal{H}(X \mid Z) = \mathcal{I}(X, Z).$$

By symmetry, we must also have $\mathcal{I}(X, Y) \leq \mathcal{I}(Y, Z)$ so that (iv) is true. ∎

Statement (iv) is often called the data-processing inequality, and has been used to infer conditional independence within gene networks [125].

## 13.2 Models on undirected graphs

### 13.2.1 Graphical representation of conditional independence

An undirected graph is a collection of vertexes and edges, in which edges link pairs of vertexes without order. Edges can therefore be identified to *subsets* of cardinality two of the set of vertexes, $V$. This yields the definition:

**Definition 13.10** *An undirected graph $G$ is a pair $G = (V, E)$ where $V$ is a finite set of vertexes and elements $e \in E$ are subsets $e = \{s, t\} \subset V$.*

Note that edges in undirected graphs are defined as sets, i.e., unordered pairs, which are delimited with braces in these notes. Later on, we will use parentheses to represent ordered pairs, $(s, t) \neq (t, s)$. We will write $s \sim_G t$, or simply $s \sim t$ to indicate that $s$ and $t$ are connected by an edge in $G$ (we also say that $s$ and $t$ are neighbors in $G$).

**Definition 13.11** *A path in an undirected graph $G = (V, E)$ is a finite sequence $(s_0, \ldots, s_N)$ of vertexes such that $s_{k-1} \sim s_k \in E$. (A sequence, $(s_0)$, of length 1 is also a path by extension.)*

*We say that $s$ and $t$ are connected by a path if either $s = t$ or there exists a path $(s_0, \ldots, s_N)$ such that $s_0 = s$ and $s_N = t$.*

*A subset $S \subset G$ is connected if any pair of elements in $S$ can be connected by a path.*

*A subset $T \subset G$ separates two other subsets $S$ and $S'$ if all paths between $S$ and $S'$ must pass in $T$. We will write $(S \perp\!\!\!\perp S' \,|\, T)$ in such a case.*

One of the goals of this chapter is to relate the notion of conditional independence within a set of variables to separation in a suitably chosen undirected graph with vertexes in one-to-one correspondence with the variables. This will also justifies the similarity of notation used for separation and conditional independence.

We have the following simple fact:

**Lemma 13.12** *Let $G = (V, E)$ be an undirected graph, and $S, S', T \subset V$. Then*

$$(S \perp\!\!\!\perp S' \,|\, T) \Rightarrow S \cap S' \subset T.$$

Indeed, if $(S \perp\!\!\!\perp S' \,|\, T)$ and $s_0 \in S \cap S'$, the path $(s_0)$ links $S$ and $S'$ and therefore must pass in $T$.

Proposition 13.5 translates into similar properties for separation:

**Proposition 13.13** *Let $(V, E)$ be an undirected graph and $S, T, U, R$ be subsets of $V$. The following properties hold*

*(i)* $(S \perp\!\!\!\perp T | U) \Leftrightarrow (T \perp\!\!\!\perp S | U)$.

*(ii)* $(S \perp\!\!\!\perp T \cup R | U) \Rightarrow (S \perp\!\!\!\perp T | U)$.

*(iii)* $(S \perp\!\!\!\perp T \cup R | U) \Rightarrow (S \perp\!\!\!\perp T | U \cup R)$.

*(iv)* $(S \perp\!\!\!\perp T \,|\, U)$ *and* $(S \perp\!\!\!\perp R \,|\, U \cup T) \Leftrightarrow (S \perp\!\!\!\perp T \cup R \,|\, U)$.

*(v)* $U \cap R = \emptyset, (S \perp\!\!\!\perp R \,|\, U \cup T)$ *and* $(S \perp\!\!\!\perp U \,|\, T \cup R) \Rightarrow (S \perp\!\!\!\perp U \cup R \,|\, T)$.

PROOF (i) is obvious, and for (ii) (and (iii)), if any path between $S$ and $T \cup R$ must pass by $U$, the same is obviously true for a path between $S$ and $T$.

For the $\Rightarrow$ part of (iv), if a path links $S$ and $T \cup R$, then it either links $S$ and $T$ and must pass through $U$ by the first assumption, or link $S$ and $R$ and therefore pass through $U$ or $T$ by the second assumption. But if the path passes through $T$, it must

also pass through $U$ before by the first assumption. In all cases, the path passes through $U$. The $\Leftarrow$ part of (iv) is obvious.

Finally, consider (v) and take a path between two distinct elements in $S$ and $U \cup R$. Consider the first time the path hits $U$ or $R$, and assumes that it hits $U$ (the other case being treated similarly by symmetry). Notice that the path cannot hit both $U$ and $R$ at the same point since $U \cap R = \emptyset$. From the assumptions, the path must hit $T \cup R$ before passing by $U$, and the intersection cannot be in $R$, so it is in $T$, which is the conclusion we wanted. ∎

To make a connection between separation in graphs and conditional independence between random variables, we consider a graph $G = (V, E)$ and a family of random variables $(X^{(s)}, s \in V)$ indexed by $V$. Each variable is assumed to take values in a set $F_s = \mathcal{R}_{X^{(s)}}$. The collection of values taken by the random variables will be called configurations, and the sets $F_s, s \in V$ are called the state spaces.

Letting $\mathbb{F}$ denote the collection $(F_s, s \in V)$, we will denote the set of such configurations as $\mathcal{F}(V, \mathbb{F})$. Then $\mathbb{F}$ is clear from the context, we will just write $\mathcal{F}(V)$. If $S \subset V$ and $x \in \mathcal{F}(V, \mathbb{F})$, the restriction of $x$ to $S$ is denoted $x^{(S)} = (x^{(s)}, s \in S)$. The set formed by those restrictions will be denoted $\mathcal{F}(S, \mathbb{F})$ (or just $\mathcal{F}(S)$).

**Remark 13.14** Some care needs to be given to the definition of the space of configurations, to avoid ambiguities when two sets $F_s$ coincide. The configuration $x = (x^{(s)}, s \in V)$ should be understood, in an emphatic way, as the collection $\hat{x} = ((s, x^{(s)}), s \in V)$, which makes explicit the fact that $x^{(s)}$ is the value observed at vertex $s$. Similarly the emphatic notation for $x^{(S)} \in \mathcal{F}(V, \mathbb{F})$ is $\hat{x}^{(S)} = ((s, x^{(s)}), s \in S)$.

In the following, we will not use the emphatic notation to avoid overly heavy expressions, but its relevance should be clear with the following simple example. Take $V = \{1, 2, 3\}$ and $F_1 = F_2 = F_3 = \{0, 1\}$. Let $x^{(1)} = 0$, $x^{(2)} = 0$ and $x^{(3)} = 1$. Then the sub-configurations $x^{(\{1,3\})}$ and $x^{(\{2,3\})}$ both corresponds to values $(0, 1)$, but we consider them as distinct. In the same spirit, $x^{(1)} = x^{(2)}$, but $x^{(\{1\})} \neq x^{(\{2\})}$ ◆

If $S, T \subset V$ with $S \cap T = \emptyset$, $x^{(S)} \in \mathcal{F}(S, \mathbb{F})$, $y^{(T)} \in \mathcal{F}(T, \mathbb{F})$, we will denote their concatenation by $x^{(S)} \wedge y^{(T)}$, which is the configuration $z = (z_s, s \in S \cup T) \in \mathcal{F}(T \cup S, \mathbb{F})$ such that $z^{(s)} = x^{(s)}$ if $s \in S$ and $z^{(s)} = y^{(s)}$ if $s \in T$.

We define a *random field* over $V$ as a random configuration $X : \Omega \to \mathcal{F}(V, \mathbb{F})$, that we will denote for short $X = (X^{(s)}, s \in V)$. If $S \subset V$, the restriction $X^{(S)}$ will also be denoted $(X^{(s)}, s \in S)$.

We can now write the definition:

**Definition 13.15** *Let $G = (V, E)$ be an undirected graph and $X = (X^{(s)}, s \in V)$ a random field over $V$. We say that $X$ is Markov (or has the Markov property) relative to $G$ (or is $G$-Markov, or is a Markov random field on $G$) if and only if, for all $S, T, U \subset V$:*

$$(S \perp\!\!\!\perp T \mid U) \Rightarrow (X^{(S)} \perp\!\!\!\perp X^{(T)} \mid X^{(U)}). \tag{13.16}$$

Letting the observation over an empty set $S$ be empty, i.e., $X_\emptyset = \emptyset$, this definition includes the statement that, if $S$ and $T$ are disconnected (i.e., there is no path between them: they are separated by the empty set), then $(X^{(S)} \perp\!\!\!\perp X^{(T)} \mid \emptyset)$: $X^{(S)}$ and $X^{(T)}$ are independent.

We will say that a probability distribution $\pi$ on $\mathcal{F}(V)$ is $G$-Markov if its associated canonical random field $X = (X^{(s)}, s \in V)$ defined on $\tilde{\Omega} = \mathcal{F}(V)$ by $X^{(s)}(x) = x^{(s)}$ is $G$-Markov.

### 13.2.2   Reduction of the Markov property

We now proceed, in a series of steps, to a simplification of definition 13.15 in order to obtain a minimal number of conditional independence statements. Note that, in its current form, definition 13.15 requires to check (13.16) for any three subsets of $V$, which provides a huge number of conditions. Fortunately, as we will see, these conditions are not independent, and checking a much smaller number of them will ensure that all of them are true.

The first step for our reduction is provided by the following lemma.

**Lemma 13.16** *Let $G = (V, E)$ be an undirected graph and $X = (X_s, s \in V)$ a set of random variables indexed by $V$. Then $X$ is $G$-Markov if and only if, for $S, T, U \subset V$,*

$$S \cap U = T \cap U = \emptyset \text{ and } (S \perp\!\!\!\perp T \mid U) \Rightarrow (X^{(S)} \perp\!\!\!\perp X^{(T)} \mid X^{(U)}). \tag{13.17}$$

Proof Assume that (13.17) is true, and take any $S, T, U$ with $(S \perp\!\!\!\perp T \mid U)$. Let $A = S \cap U$, $B = T \cap U$ and $C = A \cup B$. Partition $S$ in $S = S_1 \cup A$, $T$ in $T_1 \cup B$ and $U$ in $U_1 \cup C$. From $(S \perp\!\!\!\perp T \mid U)$, we get $(S_1 \perp\!\!\!\perp T_1 \mid U)$. Since $S_1 \cap U = T_1 \cap U = \emptyset$, this implies $(X^{(S_1)} \perp\!\!\!\perp X^{(T_1)} \mid X^{(U)})$. But this implies $((X^{(S_1)}, X^{(A)}) \perp\!\!\!\perp (X^{(T_1)}, X^{(B)}) \mid X^{(U)})$. Indeed, this property requires

$$P_X(x^{(S_1)} \wedge x^{(A)} \wedge x^{(T_1)} \wedge x^{(B)} \wedge x^{(U_1)} \wedge y^{(C)}) P^X(x^{(U_1)} \wedge y^{(C)})$$
$$= P_X(x^{(S_1)} \wedge x^{(A)} \wedge x^{(U_1)} \wedge y^{(C)}) P_X(x^{(T_1)} \wedge x^{(B)} \wedge x^{(U_1)} \wedge y^{(C)})$$

If the configurations $x^{(A)}, x^{(B)}, y^{(C)}$ are not consistent (i.e., $x^{(t)} \neq y^{(t)}$ for some $t \in C$), then both sides vanish. So we can assume $x^{(C)} = y^{(C)}$ and remove $x^{(A)}$ and $x^{(B)}$ from the expression, since they are redundant. The resulting identity is true since it exactly states that $(X^{(S_1)} \perp\!\!\!\perp X^{(T_1)} \mid X^{(U)})$.                                        ∎

Define the set of neighbors of $s \in V$ (relative to the graph $G$) as the set of $t \neq s$ such that $\{s, t\} \in E$ and denote this set by $\mathcal{V}_s$. For $S \subset V$ define also

$$\mathcal{V}_S = S^c \cap \bigcup_{s \in S} \mathcal{V}_s$$

which is the set of neighbors of all vertexes in $S$ that do not belong to $S$. (Here $S^c$ denotes the complementary set of $S$, $S^c = V \setminus S$.) Finally, let $\mathcal{W}_S$ denote the vertexes that are "remote" from $S$, $\mathcal{W}_S = (S \cup \mathcal{V}_S)^c$.

We have the following important reduction of the condition in definition 13.15.

**Proposition 13.17** *X is Markov relative to G if and only if, for any $S \subset V$,*

$$(X^{(S)} \perp\!\!\!\perp X^{(\mathcal{W}_S)} \mid X^{(\mathcal{V}_S)}). \tag{13.18}$$

This says that

$$\mathbb{P}(X^{(S)} = x^{(S)} \mid X^{(S^c)} = x^{(S^c)})$$

only depends (when defined) on variables $x^{(t)}$ for $t \in S \cup \mathcal{V}_S$.

PROOF First note that $(S \perp\!\!\!\perp \mathcal{W}_S \mid \mathcal{V}_S)$ is always true, since any path reaching $S$ from $S^c$ must pass through $\mathcal{V}_S$. This immediately proves the "only if" part of the proposition.

Consider now the "if" part. Take $S, T, U$ such that $(S \perp\!\!\!\perp T \mid U)$. We want to prove that $(X_S \perp\!\!\!\perp X_T \mid X_U)$. According to lemma 13.16, we can assume, without loss of generality, that $S \cap U = T \cap U = \emptyset$.

Define $R$ as the set of vertexes $v$ in $V$ such that there exists a path between $S$ and $v$ that does not pass in $U$. Then:

1. $S \subset R$: the path $(s)$ for $s \in S$ does not pass in $U$ since $S \cap U = \emptyset$.

2. $U \cap R = \emptyset$ by definition.

3. $\mathcal{V}_R \subset U$: assume that there exists a point $r$ in $\mathcal{V}_R$ which is not in $U$. Then $r$ has a neighbor, say $r'$ in $R$. By definition of $R$, there exists a path from $S$ to $r'$ that does not hit $U$, and this path can obviously be extended by adding $r$ at the end to obtain a path that still does not hit $U$. But this implies that $r \in R$, which contradicts the fact that $\mathcal{V}_R \cap R = \emptyset$.

4. $T \cap (R \cup \mathcal{V}_R) = \emptyset$: if $t \in T$, then $t \notin R$ from $(S \perp\!\!\!\perp T \mid U)$ and $t \notin \mathcal{V}_R$ from $T \cap U = \emptyset$.

We can then write (each decomposition being a partition, implicitly defining the sets $A$, $B$ and $C$, see Fig. 13.1) $R = S \cup A$, $U = \mathcal{V}_R \cup C$, $(R \cup \mathcal{V}_R)^c = T \cup C \cup B$, and from $(X^{(R)} \perp\!\!\!\perp X^{(\mathcal{W}_R)} \mid X^{(\mathcal{V}_R)})$, we get

$$((X^{(S)}, X^{(A)}) \perp\!\!\!\perp (X^{(T)}, X^{(C)}, X^{(B)}) \mid X^{(\mathcal{V}_R)})$$

Figure 13.1: See proof of proposition 13.17 for details

which implies

$$((X^{(S)}, X^{(A)}) \perp\!\!\!\perp (X^{(T)}, X^{(B)}) \mid X^{(U)})$$

by (CI3), which finally implies $(X^{(S)} \perp\!\!\!\perp X^{(T)} \mid X^{(U)})$ by (CI2).                    ∎

For positive probabilities, it suffices to consider singletons in proposition 13.17.

**Proposition 13.18** *If the joint distribution of $(X^{(s)}, s \in V)$ is positive and, for any $s \in V$,*

$$(X^{(s)} \perp\!\!\!\perp X^{(\mathcal{W}_s)} \mid X^{(\mathcal{V}_s)}), \tag{13.19}$$

*then X is Markov relative to G.  The converse statement is true without the positivity assumption.*

Proof  It suffices to prove that, if (13.18) is true for $S$ and $T \subset V$, with $T \cap S = \emptyset$, it is also true for $S \cup T$. The result will then follow by induction.

So, let $U = \mathcal{V}_{S \cup T}$ and $R = \mathcal{W}_{S \cup T} = V \setminus (S \cup T \cup U)$. Then, we have

$$(X^{(S)} \perp\!\!\!\perp X^{(\mathcal{W}_S)} \mid X^{(\mathcal{V}_S)}) \Rightarrow (X^{(S)} \perp\!\!\!\perp X^{(R)} \mid (X^{(U)}, X^{(T)}))$$

because $R \subset \mathcal{W}_S$ (if $s \in \mathcal{V}_S$, then it is either in $U$ or in $T$ and therefore cannot be in $R$). Similarly, $(X^{(T)} \perp\!\!\!\perp X^{(R)} \mid (X^{(U)}, X^{(S)}))$, and (CI5) (for which we need $P$ positive) now implies $((X^{(T)}, X^{(S)}) \perp\!\!\!\perp X^{(R)} \mid X^{(U)})$.                    ∎

To see that the positivity assumption is needed, consider the following example with six variables $X^{(1)}, \ldots, X^{(6)}$, and a graph linking consecutive integers and closing with

an edge between 1 and 6. Assume that $X^{(1)} = X^{(2)} = X^{(4)} = X^{(5)}$, and that $X^{(1)}, X^{(3)}$ and $X^{(6)}$ are independent. Then (13.19) is true, since, for $k = 1, 2, 4, 5$, $X^{(k)}$ is constant given its neighbors, and $X^{(3)}$ (resp. $X^{(6)}$) is independent of the rest of the variables. But $(X^{(1)}, X^{(2)})$ is not independent of $(X^{(4)}, X^{(5)})$ given the neighbors $X^{(3)}, X^{(6)}$.

Finally, another statement equivalent to proposition 13.18 is the following:

**Proposition 13.19** *If the joint distribution of* $(X^{(s)}, s \in V)$ *is positive and, for any* $s, t \in V$,

$$s \nsim_G t \Rightarrow (X^{(s)} \perp\!\!\!\perp X^{(t)} \mid X^{(V \setminus \{s,t\})}),$$

*then X is Markov relative to G. The converse statement is true without the positivity assumption.*

PROOF Fix $s \in V$ and assume that $(X^{(s)} \perp\!\!\!\perp X^{(R)} \mid X^{(V \setminus R)})$ for any $R \subset \mathcal{W}_s$ with cardinality at most $k$ (the statement is true for $k = 1$ by assumption). Consider a set $\tilde{R} \subset \mathcal{W}_s$ of cardinality $k + 1$, that we decompose into $R \cup \{t\}$ for some $t \in \tilde{R}$. We have $(X^{(s)} \perp\!\!\!\perp X^{(t)} \mid X^{(V \setminus \tilde{R})}, X_R)$ from the initial hypothesis and $(X^{(s)} \perp\!\!\!\perp X^{(R)} \mid X^{(V \setminus \tilde{R})}, X_t)$ from the induction hypothesis. Using property (CI5), this yields $(X^{(s)} \perp\!\!\!\perp X^{(\tilde{R})} \mid X^{(V \setminus \tilde{R})})$. This proves the proposition by induction. ∎

**Remark 13.20** It is obvious from the definition of a $G$-Markov process that, if $X$ is Markov for a graph $G = (V, E)$, it is automatically Markov for any richer graph, i.e., any graph $\tilde{G} = (V, \tilde{E})$ with $E \subset \tilde{E}$. This is because separation in $\tilde{G}$ implies separation in $G$. Moreover, any $X$ is $G$-Markov for the *complete graph* on $V$, for which $s \sim t$ for all $s \neq t \in V$. This is because no pair of sets can be separated in a complete graph.

Any graph with respect to which $X$ is Markov must be richer than the graph $G_X = (V, E_X)$ defined by $s \nsim_{G_X} t$ if and only $(X^{(s)} \perp\!\!\!\perp X^{(t)} \mid X^{(\{s,t\}^c)})$. This is true because, for any graph $G$ for which $X$ is Markov, we have

$$s \nsim_G t \Rightarrow (X^{(s)} \perp\!\!\!\perp X^{(t)} \mid X^{(\{s,t\}^c)}) \Rightarrow s \nsim_{G_X} t.$$

Interestingly, proposition 13.19 states that $X$ is $G_X$-Markov as soon as its joint distribution is positive. This implies that $G_X$ is the minimal graph over which $X$ is Markov in this case. ♦

### 13.2.3 Restricted graph and partial evidence

Assume that some variables $X^{(T)} = (X^{(t)}, t \in T)$ (with $T \subset V$) have been observed, with observed values $x^{(T)} = (x^{(t)}, t \in T)$. One would like to use this partial evidence to get additional information on the remaining variables, $X^{(S)}$ where $S = V \setminus T$. From the probabilistic point of view, this means computing the conditional distribution of $X^{(S)}$ given $X^{(T)} = x^{(T)}$.

One important property of $G$-Markov models is that the Markov property is essentially conserved when passing to conditional distributions. We introduce for this the following definitions.

**Definition 13.21** *If $G = (V, E)$ is an undirected graph, a subgraph of $G$ is a graph $G' = (V', E')$ with $V' \subset V$ and $E' \subset E$.*

*If $S \subset V$, the* restricted graph, $G_S$, *of $G$ to $S$ is defined by*

$$G_S = (S, E_S) \text{ with } E_S = \{e = \{s, t\} : s, t \in S \text{ and } e \in E\}. \tag{13.20}$$

We have the following proposition.

**Proposition 13.22** *Let $G = (V, E)$ be an undirected graph and $X$ be $G$-Markov. Let $S \subset V$ and $T = S^c$. Given a partial evidence $x^{(T)}$ such that $P(X^{(T)} = x^{(T)}) > 0$, $X^{(S)}$, conditionally to $X^{(T)} = x^{(T)}$, is $G_S$-Markov.*

PROOF  The proof is straightforward once it is noticed that

$$(A \perp\!\!\!\perp B \mid C)_{G_S} \Rightarrow (A \perp\!\!\!\perp B \mid C \cup T)_G \qquad\qquad\blacksquare$$

so that

$$\begin{aligned}
(A \perp\!\!\!\perp B \mid C)_{G_S} &\Rightarrow (X^{(A)} \perp\!\!\!\perp X^{(B)} \mid X^{(C)}, X^{(T)})_P \\
&\Rightarrow (X^{(A)} \perp\!\!\!\perp X^{(B)} \mid X^{(C)})_{P(\cdot \mid X^{(T)} = x^{(T)})}
\end{aligned}$$

### 13.2.4  Marginal distributions

The effect of taking marginal distributions for a $G$-Markov model is, unfortunately, not as much a mild operation as computing conditional distributions, in the sense that the conditional independence structure of the marginal distribution may be much more complex than the original one.

Let $G = (V, E)$ be an undirected graph, and let $S$ be a subset of $V$. Define the graph $G^S = (S, E^S)$ by $\{s, t\} \in E^S$ if and only if $\{s, t\} \in E$ or there exist $u, u' \in S^c$ such that $\{s, u\} \in E$, $\{t, u'\} \in E$ and $u$ and $u'$ are connected by a path in $S^c$. In other terms $E^S$ links all $s, t \in S$ that can be connected by a path, all but the extremities of which are included in $S^c$. With this notation, the following proposition holds.

**Proposition 13.23** *Let $G = (V, E)$ be an undirected graph, and $S \subset V$. Assume that $X = (X^{(s)}, s \in V)$ is a family of random variables which is $G$-Markov. Then $X^{(S)} = (x^{(s)}, s \in S)$ is $G^S$-Markov.*

PROOF It suffices to prove that, for $A, B, C \subset S$,

$$(A \perp\!\!\!\perp B \mid C)_{G^S} \Rightarrow (A \perp\!\!\!\perp B \mid C)_G.$$

So, assume that $A$ and $B$ are separated by $C$ in $G^S$. If a path connects $A$ and $B$ in $G$, we can, by definition of $E^S$, remove from this path any portion that passes in $S^c$ and obtain a valid path in $G^S$. By assumption, this path must pass in $C$, and therefore so does the original path. ∎

The graph $G^S$ can be much more complex than the restricted graph $G_S$ introduced in the previous section (note that, by definition, $G^S$ is richer than $G_S$). Take, for example, the graph that corresponds to "hidden Markov models," for which (cf. fig. 13.2)
$$V = \{1, \dots, N\} \times \{0, 1\}$$
and edges $\{s, t\} \in E$ have either $s = (k, 0)$ and $t = (l, 0)$ with $|k - l| = 1$, or $s = (k, 0)$ and $t = (k, 1)$. Let $S = \{1, \dots, N\} \times \{1\}$. Then, $G_S$ is totally disconnected ($E_S = \emptyset$), since no edge in $G$ links two elements of $S$. In contrast, any pair of elements in $S$ is connected by a path in $S^c$, so that $G^S$ is a complete graph.



Figure 13.2: In this graph, variables in the lower row are conditionally independent given the first row, while their marginal distribution requires a completely connected graph.

## 13.3   The Hammersley-Clifford theorem

The Hammersley-Clifford theorem, which will be proved in this section, gives a complete description of positive Markov processes relative to a given graph, $G$. It states that positive $G$-Markov models are associated to families of positive local interactions indexed by cliques in the graph. We now introduce each of these concepts.

### 13.3.1   Families of local interactions

**Definition 13.24** *Let $V$ be a set of vertexes and $(F_s, s \in V)$ a collection of state spaces. A family of local interactions is a collection of non-negative functions $\Phi = (\varphi_C, C \in \mathcal{C})$ indexed over some subset $\mathcal{C}$ of $\mathcal{P}(V)$, such that each $\varphi_C$ only depends on configurations*

*restricted to C (i.e., it is defined on $\mathcal{F}(C)$), with values in $[0, +\infty)$. (Recall that $\mathcal{P}(V)$ is the set of all subsets of V.)*

*Such a family has order p if no $C \in \mathcal{C}$ has cardinality larger than p. A family of local interactions of order 2 is also called a family of pair interactions.*

*Such a family is said to be consistent, if there exists an $x \in \mathcal{F}(V)$ such that*

$$\prod_{C \in \mathcal{C}} \varphi_C(x^{(C)}) \neq 0.$$

*To a consistent family of local interactions, one associates the probability distribution $\pi^\Phi$ on $\mathcal{F}(V)$ defined by*

$$\pi^\Phi(x) = \frac{1}{Z^\Phi} \prod_{C \in \mathcal{C}} \varphi_C(x^{(C)}) \tag{13.21}$$

*for all $x \in \mathcal{F}(V)$, where $Z^\Phi$ is a normalizing constant.*

Given $\mathcal{C} \subset \mathcal{P}(V)$, define the graph $G_\mathcal{C} = (V, E_\mathcal{C})$ by letting $\{s, t\} \in E_\mathcal{C}$ if and only if there exists $C \in \mathcal{C}$ such that $\{s, t\} \in C$. We then have the following proposition.

**Proposition 13.25** *Let $\Phi = (\varphi_C, C \subset \mathcal{C})$ be a consistent family of local interactions, associated to some $\mathcal{C} \subset \mathcal{P}(V)$. Then the associated distribution $\pi^\Phi$ is $G_\mathcal{C}$-Markov.*

PROOF Let $X$ be a random field associated with $\pi = \pi^\Phi$. According to proposition 13.17, we must show that, for any $S \subset V$, one has

$$(X^{(S)} \perp\!\!\!\perp X^{(\mathcal{W}_S)} \mid X^{(\mathcal{V}_S)})$$

where $\mathcal{V}_S$ is the set of neighbors of $S$ in $G_\mathcal{C}$ and $\mathcal{W}_S = V \setminus (\mathcal{V}_S \cup S)$. Define the set $U_S$ by

$$U_S = \bigcup_{C \in \mathcal{C}, S \cap C \neq \emptyset} C$$

so that $\mathcal{V}_S = U_S \setminus S$ and $\mathcal{W}_S = V \setminus U_S$. To prove conditional independence, we need to prove that, for any $x \in F$:

$$\pi(x)\pi_{\mathcal{V}_S}(x^{(\mathcal{V}_S)}) = \pi_{U_S}(x^{(U_S)})\pi_{V \setminus S}(x^{(V \setminus S)}) \tag{13.22}$$

(where we denote $\pi_A$ the marginal distribution of $\pi$ on $\mathcal{F}(A)$.)

From the definition of $\pi$, we have

$$\begin{aligned}
\pi(x) &= \frac{1}{Z} \prod_{C \in \mathcal{C}} \varphi_C(x^{(C)}) \\
&= \frac{1}{Z} \prod_{C : C \cap S \neq \emptyset} \varphi_C(x^{(C)}) \prod_{C : C \cap S = \emptyset} \varphi_C(x^{(C)}).
\end{aligned}$$

The first term in the last product only depends on $x^{(U_S)}$, and the second one only on $x^{(V \setminus S)}$. Introduce the notation

$$
\begin{cases}
\mu_1(x^{(\mathcal{V}_S)}) = \displaystyle\sum_{y^{(U_S)}: y^{(\mathcal{V}_S)} = x^{(\mathcal{V}_S)}} \prod_{C: C \cap S \neq \emptyset} \varphi_C(x^{(C)}) \\[2ex]
\mu_2(x^{(\mathcal{V}_S)}) = \displaystyle\sum_{y^{(V \setminus S)}: y^{(\mathcal{V}_S)} = x^{(\mathcal{V}_S)}} \prod_{C: C \cap S = \emptyset} \varphi_C(x^{(C)}).
\end{cases}
$$

With this notation, we have:

$$
\begin{cases}
\pi_{U_S}(x^{(U_S)}) = (\mu_2(x^{(\mathcal{V}_S)})/Z) \displaystyle\prod_{C: C \cap S \neq \emptyset} \varphi_C(x^{(C)}) \\[2ex]
\pi_{V \setminus S}(x^{(V \setminus S)}) = (\mu_1(x^{(\mathcal{V}_S)})/Z) \displaystyle\prod_{C: C \cap S = \emptyset} \varphi_C(x^{(C)}) \\[2ex]
\pi_{\mathcal{V}_S}(x^{(\mathcal{V}_S)}) = \mu_1(x^{(\mathcal{V}_S)}) \mu_2(x^{(\mathcal{V}_S)})/Z
\end{cases}
$$

from which (13.22) can be easily obtained. ∎

We now discuss conditional distributions and marginals for processes associated with local interactions. If $T \subset V$, we let $\pi_T = \pi_T^\Phi$ denote the marginal distribution of $\pi$ on $T$.

We start with a discussion of conditionals. Let $\pi$ be associated with $\Phi$, and let $S \subset V$ and $T = V \setminus S$. Assume that a configuration $y^{(T)}$ is given, such that $\pi_T(y^{(T)}) > 0$, and consider the conditional distribution

$$
\pi_{S|T}(x^{(S)} \mid y^{(T)}) = \pi(x^{(S)} \wedge y^{(T)})/\pi_T(y^{(T)}). \tag{13.23}
$$

We have the following proposition.

**Proposition 13.26** *With the notation above, $\pi_{S|T}(\cdot \mid y^{(T)})$ is associated to the family of local interactions $\Phi_{|y_T} = (\varphi_{\tilde{C}|y^{(T)}}, \tilde{C} \in \mathcal{C}_S)$ with*

$$
\mathcal{C}_S = \left\{ \tilde{C} : \tilde{C} \subset S, \exists C \in \mathcal{C} : \tilde{C} = C \cap S \right\}
$$

*and*

$$
\varphi_{\tilde{C}|y^{(T)}}(x^{(\tilde{C})}) = \prod_{C \in \mathcal{C}: C \cap S = \tilde{C}} \varphi_C(x^{(\tilde{C})} \wedge y^{(C \cap T)}).
$$

PROOF From (13.23) and the definition of $\pi$, it is easy to sees that

$$
\pi_{S|T}(x^{(S)} \mid y^{(T)}) = \frac{1}{Z(y^{(T)})} \prod_{C: C \cap S \neq \emptyset} \varphi_C(x^{(C \cap S)} \wedge y^{(C \cap T)}),
$$

where $Z(y^{(T)})$ is a constant that only depends on $y^{(T)}$. The fact that $\pi_{S|T}(\cdot \mid y^{(T)})$ is associated to $\Phi_{|y^{(T)}}$ is then obtained by reorganizing the product over distinct $S \cap C$'s. ∎

This result, combined with proposition 13.25, is consistent with proposition 13.22, in the sense that the restriction to $G_{\mathcal{C}}$ to $S$ coincides with the graph $G_{\mathcal{C}_S}$. The easy proof is left to the reader.

We now consider marginals, and more specifically marginals when only one node is removed, which provides the basis for "node elimination."

**Proposition 13.27** *Let $\pi$ be associated to $\Phi = (\varphi_C, C \in \mathcal{C})$ as above. Let $t \in V$ and $S = V \setminus \{t\}$. Define $\mathcal{C}_t \in \mathcal{P}(V)$ as the set*

$$\mathcal{C}_t = \{C \in \mathcal{C} : t \notin C\} \cup \{\tilde{C}_t\}$$

*with*

$$\tilde{C}_t = \bigcup_{C \in \mathcal{C}: t \in C} C \setminus \{t\}.$$

*Define a family of local interactions $\Phi_t = (\tilde{\varphi}_{\tilde{C}}, \tilde{C} \in \mathcal{C}_t)$ by $\tilde{\varphi}_{\tilde{C}} = \varphi_{\tilde{C}}$ if $\tilde{C} \neq \tilde{C}_t$ and:*

- *If $\tilde{C}_t \notin \mathcal{C}$:*

$$\tilde{\varphi}_{\tilde{C}_t}(x^{(\tilde{C}_t)}) = \sum_{y^{(t)} \in F_t} \prod_{C \in \mathcal{C}, t \in C} \varphi_C(x^{(\tilde{C}_t)} \wedge y^{(t)}).$$

- *If $\tilde{C}_t \in \mathcal{C}$:*

$$\tilde{\varphi}_{\tilde{C}_t}(x^{(\tilde{C}_t)}) = \varphi_{C_t}(x^{(\tilde{C}_t)}) \sum_{y^{(t)} \in F_t} \prod_{C \in \mathcal{C}, t \in C} \varphi_C(x^{(C_t)} \wedge y^{(t)})$$

*Then the marginal, $\pi_S$, of $\pi$ over $S$ is the distribution associated to $\Phi_t$.*

The proof is almost straightforward by summing over possible values of $y_t$ in the expression of $\pi$ and left to the reader.

### 13.3.2  Characterization of positive $G$-Markov processes

Using families of local interactions is a typical way to build graphical models in applications. The previous section describes a graph with respect to which the obtained process is Markov. Conversely, given a graph $G$, the Hammersley-Clifford theorems states that families of local interactions over the cliques of $G$ are the only ways to build positive graphical models, which reinforces the importance of this construction. We now pass to the statement and proof of this theorem, starting with the following definition.

**Definition 13.28** *Let $G = (V, E)$ be an undirected graph. A clique in $G$ is a nonempty subset $C \subset V$ such that $s \sim_G t$ whenever $s, t \in C$, $s \neq t$. (In particular, subsets of cardinality one are always cliques.) Cliques therefore form complete subgraphs of $G$.*

*The set of cliques of a graph $G$ will be denoted $\mathcal{C}_G$.*

*A clique that cannot be strictly included in any other clique is called a maximal clique, and their set denoted $\mathcal{C}_G^*$.*

*(Note that some authors call cliques what we refer to as maximal cliques.)*

Given $G = (V, E)$, consider a random field $X = (X^{(s)}, s \in V)$. We assume that $X^{(s)}$ takes values in a finite set $F_s$ with $\mathbb{P}(X^{(s)} = a) > 0$ for any $a \in F_s$ (this is no loss of generality since one can always restrict $F_s$ to such $a$'s). If $S \subset V$, we denote as before $\mathcal{F}(S)$ the set of restrictions of configurations to $S$. With this notation, $X$ is positive, according to definition 13.4, if and only if $P(X = x) > 0$ for all $x \in \mathcal{F}(V)$. We will let $\pi = P_X$ be the probability distribution of $X$, so that $\pi(x) = \mathbb{P}(X = x)$ and use as above the notation: for $S, T \subset V$

$$\begin{cases} \pi_S(x^{(S)}) = \mathbb{P}(X^{(S)} = x^{(S)}) \\ \pi_{S|T}(x^{(S)} \mid x^{(T)}) = \mathbb{P}(X^{(S)} = x^{(s)} \mid X^{(T)} = x^{(T)}). \end{cases} \tag{13.24}$$

(For the first notation, we will simply write $\pi$ if $S = V$.)

We will also need to fix a reference, or "zero," configuration in $\mathcal{F}(V)$ that we will denote $\mathbb{0} = (\mathbb{0}^{(s)}, s \in V)$, with $\mathbb{0}^{(s)} \in F_s$ for all $s$. We can choose it arbitrarily. Given this, we have the theorem:

**Theorem 13.29 (Hammersley-Clifford)** *With the previous notation, $X$ is a positive $G$-Markov process if and only if its distribution, $\pi$, is associated to a family of local interactions $\Phi = (\varphi_C, C \subset \mathcal{C}_G)$ such that $\varphi_C(x^{(C)}) > 0$ for all $x^{(C)} \in \mathcal{F}(C)$.*

*Moreover, $\Phi$ is uniquely characterized by the additional constraint: $\varphi_C(x^{(C)}) = 1$ as soon as there exists $s \in C$ such that $x^{(s)} = \mathbb{0}^{(s)}$.*

Letting $\lambda_C = -\log \varphi_C$, we get an equivalent formulation of the theorem in terms of potentials, where a potential is defined as a family of functions

$$\Lambda = (\lambda_C, C \in \mathcal{C})$$

indexed by a subset $\mathcal{C}$ of $\mathcal{P}(V)$, such that $\lambda_C$ only depends on $x^{(C)}$. The distribution associated to $\Lambda$ is

$$\pi(x) = \frac{1}{Z_\Lambda} \exp\left(-\sum_{C \in \mathcal{C}} \lambda_C(x^{(C)})\right). \tag{13.25}$$

With this terminology, we trivially have an equivalent formulation:

**Theorem 13.30** *X is a positive G-Markov process if and only if its distribution, $\pi$, is associated to a potential $\Lambda = (\lambda_C, C \subset \mathcal{C}_G)$.*

*Moreover, $\Lambda$ is uniquely characterized by the additional constraint: $\lambda_C(x^{(C)}) = 0$ as soon as there exists $s \in C$ such that $x^{(s)} = \mathbb{0}^{(s)}$.*

We now prove this theorem.

PROOF Let us start with the "if" part. If $\pi$ is associated to a potential over $\mathcal{C}_G$, we have already proved that $\pi$ is $G_{\mathcal{C}_G}$-Markov, so that it suffices to prove that $G_{\mathcal{C}_G} = G$, which is almost obvious: If $s \sim_G t$, then $\{s, t\} \in \mathcal{C}_G$ and $s \sim_{G_{\mathcal{C}_G}} t$ by definition of $G_{\mathcal{C}_G}$. Conversely, if $s \sim_{G_{\mathcal{C}_G}} t$, there exists $C \in \mathcal{C}_G$ such that $\{s, t\} \subset C$, which implies that $s \sim_G t$, by definition of a clique.

We now prove the "only if" part, which relies on a combinatorial lemma, which is one of Möbius's inversion formulas.

**Lemma 13.31** *Let A be a finite set and $f : \mathcal{P}(A) \to \mathbb{R}$, $B \mapsto f_B$. Then, there is a unique function $\lambda : \mathcal{P}(A) \to \mathbb{R}$ such that*

$$\forall B \subset A, \ f_B = \sum_{C \subset B} \lambda_C, \tag{13.26}$$

*and $\lambda$ is given by*

$$\lambda_C = \sum_{B \subset C} (-1)^{|C| - |B|} f_B. \tag{13.27}$$

To prove the lemma, first notice that the space $F$ of functions $f : \mathcal{P}(A) \to \mathbb{R}$ is a vector space of dimension $2^{|A|}$ and that the transformation $\varphi : \lambda \mapsto f$ with $f_B = \sum_{C \subset B} \lambda_C$ is linear. It therefore suffices to prove that, given any $f$, the function $\lambda$ given in (13.27) satisfies $\varphi(\lambda) = f$, since this proves that $\varphi$ is onto from $F$ to $F$ and therefore necessarily one to one.

So consider $f$ and $\lambda$ given by (13.27). Then

$$\varphi(\lambda)(B) = \sum_{C \subset B} \lambda_C = \sum_{C \subset B} \sum_{\tilde{B} \subset C} (-1)^{|C| - |\tilde{B}|} f_{\tilde{B}} = \sum_{\tilde{B} \subset B} \left( \sum_{C \supset \tilde{B}, C \subset B} (-1)^{|C| - |\tilde{B}|} \right) f_{\tilde{B}} = f_B$$

The last identity comes from the fact that, for any finite set $\tilde{B} \subset B, \tilde{B} \neq B$, we have

$$\sum_{C \supset \tilde{B}, C \subset B} (-1)^{|C| - |\tilde{B}|} = 0$$

(for $\tilde{B} = B$, the sum is obviously equal to 1). Indeed, if $s \in B$, $s \notin \tilde{B}$, we have

$$\sum_{C \supset \tilde{B}, C \subset B} (-1)^{|C|-|\tilde{B}|} = \sum_{C \supset \tilde{B}, C \subset B, s \in C} (-1)^{|C|-|\tilde{B}|} + \sum_{C \supset \tilde{B}, C \subset B, s \notin C} (-1)^{|C|-|\tilde{B}|}$$

$$= \sum_{C \supset \tilde{B}, C \subset B, s \notin C} ((-1)^{|C \cup \{s\}|-|\tilde{B}|} + (-1)^{|C|-|\tilde{B}|})$$

$$= 0.$$

So the lemma is proved. We now proceed to proving the existence and uniqueness statements in theorem 13.30. Assume that $X$ is $G$-Markov and positive. Fix $x \in \mathcal{F}(V)$ and consider the function, defined on $\mathcal{P}(V)$ by

$$f_B(x^{(B)}) = -\log \frac{\pi(x^{(B)} \wedge \mathbb{0}^{(B^c)})}{\pi(\mathbb{0})}.$$

Then, letting

$$\lambda_C(x^{(C)}) = \sum_{B \subset C} (-1)^{|C|-|B|} f_B(x^{(B)}),$$

we have $f_B(x^{(B)}) = \sum_{C \subset B} \lambda_C(x^{(C)})$. In particular, for $B = V$, this gives

$$\pi(x) = \frac{1}{Z} \exp\left(-\sum_{C \subset V} \lambda_C(x^{(C)})\right)$$

with $Z = P(\mathbb{0})$. We now prove that $\lambda_C(x^{(C)}) = 0$ if $x^{(s)} = \mathbb{0}^{(s)}$ for some $s \in V$ or if $C \notin \mathcal{C}_G$. This will prove (13.25) and the existence statement in theorem 13.30.

So, assume $x^{(s)} = \mathbb{0}^{(s)}$. Then, for any $B$ such that $s \notin B$, we have $f_B(x^{(B)}) = f_{\{s\} \cup B}(x^{(\{s\} \cup B)})$. Now take $C$ with $s \in C$. We have

$$\lambda_C(x^{(C)}) = \sum_{B \subset C, s \in B} (-1)^{|C|-|B|} f_B(x^{(B)}) + \sum_{B \subset C, s \notin B} (-1)^{|C|-|B|} f_B(x^{(B)})$$

$$= \sum_{B \subset C, s \notin B} (-1)^{|C|-|B \cup \{s\}|} f_{B \cup \{s\}}(x^{(B \cup \{s\})}) + \sum_{B \subset C, s \notin B} (-1)^{|C|-|B|} f_B(x^{(B)})$$

$$= \sum_{B \subset C, s \notin B} ((-1)^{|C|-|B \cup \{s\}|} + (-1)^{|C|-|B|}) f_B(x^{(B)})$$

$$= \mathbb{0}.$$

Now assume that $C$ is not a clique, and let $s \neq t \in C$ such that $s \nsim t$. We can write, using decompositions similar to the above,

$$\lambda_C(x^{(C)}) = \sum_{B \subset C \setminus \{s,t\}} (-1)^{|C|-|B|} \left(f_{B \cup \{s,t\}}(x^{(B \cup \{s,t\})}) - f_{B \cup \{s\}}(x^{(B \cup \{s\})}) - f_{B \cup \{t\}}(x^{(B \cup \{t\})}) + f_B(x^{(B)})\right).$$

But, for $B \subset C \setminus \{s, t\}$, we have

$$
\begin{aligned}
f_{B \cup \{s,t\}}(x^{(B \cup \{s,t\})}) - f_{B \cup \{s\}}(x^{(B \cup \{s\})}) &= -\log \frac{\pi(x^{(B \cup \{s,t\})} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})})}{\pi(x^{(B \cup \{s\})} \wedge \mathbb{0}^{(B^c \setminus \{s\})})} \\
&= \log \frac{\pi_t(x^{(t)} \mid x^{(B \cup \{s\})} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})})}{\pi_t(\mathbb{0}^{(t)} \mid x^{(B \cup \{s\})} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})})}
\end{aligned}
$$

and

$$
\begin{aligned}
f_{B \cup \{t\}}(x^{(B \cup \{t\})}) - f_B(x^{(B)}) &= -\log \frac{\pi(x^{(B \cup \{t\})} \wedge \mathbb{0}^{(B^c \setminus \{t\})})}{\pi(x^{(B)} \wedge \mathbb{0}^{(B^c)})} \\
&= \log \frac{\pi_t(x^{(t)} \mid x^{(B)} \wedge \mathbb{0}^{(B^c \setminus \{t\})})}{\pi_t(\mathbb{0}^{(t)} \mid x^{(B)} \wedge \mathbb{0}^{(B^c \setminus \{t\})})}.
\end{aligned}
$$

So, we can write

$$
\lambda_C(x^{(C)}) = \sum_{B \subset C \setminus \{s,t\}} (-1)^{|C| - |B|} \log \frac{\pi_t(x^{(t)} \mid x^{(B \cup \{s\})} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})}) \pi_t(\mathbb{0}^{(t)} \mid x^{(B)} \wedge \mathbb{0}^{(B^c \setminus \{t\})})}{\pi_t(\mathbb{0}^{(t)} \mid x_{B \cup \{s\}} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})}) \pi_t(x^{(t)} \mid x^{(B)} \wedge \mathbb{0}^{(B^c \setminus \{t\})})}
$$

which vanishes, because

$$
\pi_t(x^{(t)} \mid x^{(B \cup \{s\})} \wedge \mathbb{0}^{(B^c \setminus \{s,t\})}) = \pi_t(x^{(t)} \mid x^{(B)} \wedge \mathbb{0}^{(B^c \setminus \{t\})})
$$

when $s \not\sim t$.

To prove uniqueness, note that, for any zero-normalized $\Lambda$ satisfying (13.25), we must have $\pi(\mathbb{0}) = 1/Z$ and therefore, for any $x$,

$$
-\log \frac{\pi(x^{(B)} \wedge \mathbb{0}^{(B^c)})}{\pi(0)} = \sum_{C \subset B} \lambda_C(x^{(C)})
$$

(extending $\Lambda$ so that $\lambda_C = 0$ for $C \notin \mathcal{C}_G$).  But, from lemma 13.31, this uniquely defines $\Lambda$.  ∎

The exponential form of the distribution in the Hammersley-Clifford theorem is related to what is called a Gibbs distribution in statistical mechanics. More precisely:

**Definition 13.32** *Let $\mathcal{F}$ be a finite set and $W : \mathcal{F} \to \mathbb{R}$ be a scalar function. The Gibbs distribution with energy $W$ at temperature $T > 0$ is defined by*

$$
\pi(x) = \frac{1}{Z_T} e^{-\frac{W(x)}{T}}, \ x \in \mathcal{F}
$$

*The normalizing constant $Z_T = \sum_{y \in \mathcal{F}} \exp(-W(y)/T)$ is called the partition function.*

*If $\Lambda = (\lambda_C, C \subset V)$ is a potential then its associated energy is*

$$
W(x) = \sum_{C \subset V} \lambda_C(x^{(C)}).
$$

So the Hammersley-Clifford theorem implies that any positive $G$-Markov model is associated to a unique zero-normalized potential defined over the cliques of $G$. This representation can also be used to provide an alternate proof of proposition 13.19, which is left to the reader. Finally, one can restate proposition 13.26 in terms of potentials, yielding:

**Proposition 13.33** *Let $P$ be a Gibbs distribution associated with a zero-normalized potential $\lambda = (\lambda_C, C \subset V)$. Let $S \subset V$ and $T = S^c$. Then the conditional distribution of $X^{(S)}$ given $X^{(T)} = x^{(T)}$ is the Gibbs distribution associated with the zero-normalized potential $\tilde{\lambda} = (\tilde{\lambda}_C, C \subset S)$ where*

$$\tilde{\lambda}_C(y^{(C)}) = \sum_{C' \subset V, C' \cap S = C} \lambda_{C'}(y^{(C)} \wedge x^{(T \cap C')}).$$

## 13.4  Models on acyclic graphs

### 13.4.1  Finite Markov chains

We now review a few important examples of Markov processes $X$ associated to specific graphs $G = (V, E)$. We will always denote by $F_s$ the space in which $X^{(s)}$ takes his values, for $s \in V$.

The simplest example of $G$-Markov process (for any graph $G$) is the case when $X = (X^{(s)}, s \in V)$ is a collection of independent random variables. In this case, we can take $G_X = (V, \emptyset)$, the totally disconnected graph on $V$. Another simple fact is that, as already remarked, any $X$ is Markov for the complete graph $(V, \mathcal{P}_2(V))$ where $\mathcal{P}_2(V)$ contains all subsets of $V$ with cardinality 2.

Beyond these trivial (but nonetheless important) cases, the simplest graph-Markov processes are those associated with linear graphs, providing finite Markov chains. For this, we let $V$ be a finite ordered set, say,

$$V = \{0, \dots, N\}.$$

We say that $X$ is a finite Markov chain if, for any $k = 1, \dots, N$

$$(X^{(k)} \perp\!\!\!\perp (X^{(0)}, \dots, X^{(k-2)}) \mid X^{(k-1)}).$$

So we have the identity

$$\mathbb{P}(X^{(0)} = x^{(0)}, \dots, X^{(k)} = x^{(k)}) P(X^{(k-1)} = x^{(k-1)})$$
$$= \mathbb{P}(X^{(0)} = x^{(0)}, \dots, X^{(k-1)} = x^{(k-1)}) \mathbb{P}(X^{(k-1)} = x^{(k-1)}, X^{(k)} = x^{(k)}).$$

The distribution of a Markov chain is therefore fully specified by $\mathbb{P}(X^{(0)} = x^{(0)}), x_0 \in F_0$ (the initial distribution) and the conditional probabilities

$$p_k(x^{(k-1)}, x^{(k)}) = \mathbb{P}(X^{(k)} = x^{(k)} \mid X^{(k-1)} = x^{(k-1)}) \tag{13.28}$$

(with an arbitrary choice when $\mathbb{P}(X^{(k-1)} = x^{(k-1)}) = 0$). Indeed, assume that $\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k-1)} = x^{(k-1)})$ is known (for all $x^{(0)}, \ldots, x^{(k-1)}$). Then, either:

(i)  $\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k-1)} = x^{(k-1)}) = 0$, in which case

$$\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k)} = x^{(k)}) = 0$$

for any $x^{(k)}$, or:

(ii)  $\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k-1)} = x^{(k-1)}) > 0$, in which case, necessarily, $\mathbb{P}(X^{(k-1)} = x^{(k-1)}) > 0$, and

$$\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k)} = x^{(k)}) = p_k(x^{(k-1)}, x^{(k)})\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(k-1)} = x^{(k-1)}).$$

Note that $p_k$ in (13.28) is a transition probability (according to definition 12.2) between $F_{k-1}$ and $F_k$.

We have the following identification of a finite Markov chain with a graph-Markov process:

**Proposition 13.34** *Let $X = (X^{(0)}, \ldots, X^{(N)})$ be a finite Markov chain, such that $X$ is positive. Then $X$ is G-Markov for the linear graph $G = (V, E)$ with*

$$
\begin{aligned}
V &= \{1, \ldots, N\} \\
E &= \{\{1, 2\}, \ldots, \{N-1, N\}\}.
\end{aligned}
$$

*The converse is true without the positivity assumption: a G-Markov process for the graph above is always a finite Markov chain.*

PROOF  We prove the direct statement (the converse one being obvious). Let $s$ and $t$ be nonconsecutive distinct integers, with, say, $s < t$. From the Markov chain assumption, we have

$$(X^{(t)} \perp\!\!\!\perp (X^{(s)}, X^{(\{1, t-2\} \setminus \{s, \})}) \mid X^{(t-1)}),$$

which, using (CI3), yields $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid X^{(\{1, \ldots, t-1\} \setminus \{s\})})$. Define $Y^{(u)} = X^{(\{1, \ldots, u\} \setminus \{s, t\})}$: what we have proved is $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid Y^{(t)})$.

We now proceed by induction and assume that $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid Y^{(u)})$ for some $u \geq t$. Then, we have $(X^{(u+1)} \perp\!\!\!\perp (X^{(s)}, X^{(t)}, Y^{(u-1)}) \mid X^{(u)})$, which implies (from (CI3)) $(X^{(u+1)} \perp\!\!\!\perp X^{(t)} \mid X^{(s)}, Y^{(u)})$. Applying (CI4) to $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid Y^{(u)})$ and $(X^{(t)} \perp\!\!\!\perp X^{(u+1)} \mid X^{(s)}, Y^{(u)})$, we obtain $(X^{(t)} \perp\!\!\!\perp (X^{(s)}, X^{(u+1)}) \mid Y^{(u)})$ and finally, $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid Y^{(u+1)})$. By induction, this gives $(X^{(t)} \perp\!\!\!\perp X^{(s)} \mid Y^{(N)})$ and therefore proposition 13.19 now implies that $X$ is $G$-Markov.

(The proposition can also be proved as a consequence of the decomposition

$$\mathbb{P}(X^{(0)} = x^{(0)}, \ldots, X^{(N)} = x^{(N)}) = \mathbb{P}(X^{(0)} = x^{(0)})p_1(x^{(0)}, x^{(1)}) \ldots p_N(x^{(N-1)}, x^{(N)}).) \qquad \blacksquare$$

### 13.4.2 Undirected acyclic graph models and trees

The situation with acyclic graphs is only slightly more complex than with linear graphs, but will require a few new definitions, including those of directed graphs and trees.

The difference between directed and undirected graphs is that the edges of the former are ordered pairs, namely:

**Definition 13.35** *A (finite) directed graph $G$ is a pair $G = (V, E)$ where $V$ is a finite set of vertexes and $E$ is a subset of*

$$V \times V \setminus \{(s, s), s \in V\},$$

*which satisfies, in addition,*

$$(s, t) \in E \Rightarrow (t, s) \notin E.$$

So, for directed graphs, edges $(s, t)$ and $(t, s)$ have different meanings, and we allow at most one of them in $E$. We say that the edge $e = (s, t)$ stems from $s$ and points to $t$. The *parents* of a vertex $s$ are the vertexes $t$ such that $(t, s) \in E$, and its children are the vertexes $t$ such that $(s, t) \in E$. We will also use the notation $s \rightarrow_G t$ to indicate that $(s, t) \in E$ (compare to $s \sim_G t$ for undirected graphs).

**Definition 13.36** *A path in a directed graph $G = (V, E)$ is a sequence $(s_0, \ldots, s_N)$ such that, for all $k = 1, \ldots, N$, $s_k \rightarrow_G s_{k+1}$ (this includes the "trivial", one-vertex, paths $(s_0)$). (The definition was the same for undirected graph, replacing $s_k \rightarrow_G s_{k+1}$ by $s_k \sim_G s_{k+1}$.) For both directed and undirected cases, one says that a path is closed if $s_0 = s_N$.*

*In an undirected graph, a path is folded if it can be written as $(s_0, \ldots, s_{N-1}, s_N, s_{N-1}, \ldots, s_0)$.*

*If $G = (V, E)$ is directed, one says that $t \in V$ is a descendant of $s \in V$ (or that $s$ is an ancestor of $t$) if there exists a path starting at $s$ and ending at $t$. In particular, every vertex is both a descendant and an ancestor of itself.*

We finally define acyclic graphs.

**Definition 13.37** *A loop in a directed (resp. an undirected) graph $G$ is a path $(s_0, s_1, \ldots, s_N)$, with $N \geq 3$, such that $s_N = s_0$, which passes only once through $s_0, \ldots, s_{N-1}$ (no self-intersection except at the end).*

*A (directed or undirected) graph $G$ is acyclic if it contains no loop.*

The following property will be useful.

**Proposition 13.38** *In a directed graph, any non-trivial closed path contains a loop (i.e., one can delete vertexes from it to finally obtain a loop.)*

*In an undirected graph, any non-trivial closed path which is not a union of folded paths contains a loop.*

PROOF  Take $\gamma = (s_0, s_1, \ldots, s_N)$ with $s_N = s_0$. The path being non-trivial means $N > 1$.

First take the case of a directed graph.  Clearly, $N \geq 3$ since a two-vertex path cannot be closed in an directed graph. Consider the first occurrence of a repetition, i.e., the first index for which

$$s_j \in \{s_0, \ldots, s_{j-1}\}.$$

Then there is a unique $j' \in \{0, \ldots, j-1\}$ such that $s_{j'} = s_j$, and the path $(s_{j'}, \ldots, s_{j-1})$ must be a loop (any repetition in the sequence would contradict the fact that $j$ was the first occurrence. This proves the result in the directed case.

Consider now an undirected graph.  We can recursively remove all folded subpaths, by keeping everything but their initial point, since each such operation still provide a path at the end.  Assume that this is done, still denoting the remaining path $(s_0, s_1, \ldots, s_N)$, which therefore has no folded subpath.  We must have $N \geq 3$ since $N = 1$ implies that the original path was a union of folded paths, and $N = 2$ provides a folded path. Let, $0 \leq j' < j$ be as in the directed case. Note that one must have $j' < j - 2$, since $j' = j - 1$ would imply an edge between $j$ and itself and $j' = j - 2$ induces a folded subpath. But this implies that $(s_{j'}, \ldots, s_{j-1})$ is a loop.  ∎

Directed acyclic graphs (DAG) will be important for us, because they are associated with Bayesian networks that we will discuss later.  For now, we are interested with undirected acyclic graphs and their relation to trees, which form a subclass of directed acyclic graphs, defined as follows.

**Definition 13.39** *A forest is a directed acyclic graph with the additional requirement that each of its vertexes has at most one parent.*

*A root in a forest is a vertex that has no parent. A forest with a single root is called a tree.*

It is clear that a forest has at least one root, since one could otherwise describe a nontrivial loop by starting from a any vertex and passing to its parent until the sequence self-intersects (which must happen since $V$ is finite). We will use the following definition.

**Definition 13.40** *If $G = (V, E)$ is a directed graph, its flattened graph, denoted $G^\flat = (V, E^\flat)$ is the undirected graph obtained by forgetting the edge ordering, namely*

$$\{s, t\} \in E^\flat \Leftrightarrow (s, t) \in E \text{ or } (t, s) \in E.$$

The following proposition relates forests and undirected acyclic graphs.

**Proposition 13.41** *If $G$ is a forest, then $G^\flat$ is an undirected acyclic graph.*

*Conversely, if $G$ is an undirected acyclic graph, there exists a forest $\tilde{G}$ such that $\tilde{G}^\flat = G$.*

Proof Let $G = (V, E)$ be a forest and, in order to reach a contradiction, assume that $G^\flat$ has a loop, $s_0, \ldots, s_{N-1}, s_N = s_0$. Assume that $(s_0, s_1) \in E$; then, also $(s_1, s_2) \in E$ (otherwise $s_1$ would have two parents), and this propagates to all $(s_k, s_{k+1})$ for $k = 0, \ldots, N - 1$. But, since $s_N = s_0$, this provides a loop in $G$ which is not possible. This proves thet $G^\flat$ has no loop since the case $(s_1, s_0) \in E$ is treated similarly.

Now, let $G$ be an undirected acyclic graph. Fix a vertex $s \in V$ and consider the following procedure, in which we recursively define sets $S_k$ of processed vertexes, and $\tilde{E}_k$ of oriented edges, $k \geq 0$, initialized with $S_0 = \{s\}$ and $\tilde{E}_0 = \emptyset$.

– At step $k$ of the procedure, assume that vertexes in $S_k$ have been processed and edges in $\tilde{E}_k$ have been oriented so that $(S_k, \tilde{E}_k)$ is a forest, and that $\tilde{E}_k^\flat$ is the set of edges $\{s, t\} \in E$ such that $s, t \in S_k$ (so, oriented edges at step $k$ can only involve processed vertexes).

– If $S_k = V$: stop, the proposition is proved.

– Otherwise, apply the following construction. Let $F_k$ be the set of edges in $E$ that contain exactly one element of $S_k$.

(1) If $F_k = \emptyset$, take any $s \in V \setminus S_k$ as a new root and let $S_{k+1} = S_k \cup \{s\}$, $\tilde{E}_{k+1} = \tilde{E}_k$.
(2) Otherwise, add to $\tilde{E}_k$ the oriented edges $(s, t)$ such that $s \in S_k$ and $\{s, t\} \in F_k$, yielding $\tilde{E}_{k+1}$, and add to $S_k$ the corresponding children ($t$'s) yielding $S_{k+1}$.

We need to justify the fact that $\tilde{G}_{k+1} = (S_{k+1}, \tilde{E}_{k+1})$ above is still a forest. This is obvious after Case (1), so consider Case (2). First $\tilde{G}_{k+1}$ is acyclic, since any oriented loop is a fortiori an unoriented loop and $G$ is acyclic. So we need to prove that no vertex in $S_{k+1}$ has two parents. Since we did not add any parent to the vertexes in $S_k$ and, by assumption, $(S_k, \tilde{E}_k)$ is a forest, the only possibility for a vertex to have two parents in $S_{k+1}$ is the existence of $t$ such that there exists $s, s' \in S_k$ with $\{s, t\}$ and $\{s', t\}$ in $E$. But, since $s$ and $s'$ have unaccounted edges containing them, they cannot have been introduced in $S_k$ before the previously introduced root has been added, so they are both connected to this root: but the two connections to $t$ would create a loop in $G$ which is impossible.

So the procedure carries on, and must end with $S_k = V$ at some point since we keep adding points to $S_k$ at each step. ∎

Note that the previous proof shows there is more than one possible orientation of a connected undirected tree into a tree is not unique, although uniquely specified

once a root is chosen. The proof is constructive, and provides an algorithm building a forest from an undirected acyclic graph.

We now define graphical models supported by trees, which constitute our first Markov models associated with directed graphs. Define the depth of a vertex in a tree $G = (V, E)$ to be the number of edges in the unique path that links it to the root. We will denote by $G_d$ the set of vertexes in $G$ that are at depth $d$, so that $G_0$ contains only the root, $G_1$ the children of the root and so on. Using this, we have the definition:

**Definition 13.42** *Let $G = (V, E)$ be a tree. A process $X = (X^{(s)}, s \in V)$ is $G$-Markov if and only, for each $d \geq 1$, and for each $s \in G_d$, we have*

$$(X^{(s)} \perp\!\!\!\perp (X^{(G_d \setminus \{s\})}, X^{(G_q \setminus \{pa(s)\})}, q < d) \mid X^{(pa(s))}) \tag{13.29}$$

*where $pa(s)$ is the parent of $s$.*

So, conditional to its parent, $X^{(s)}$ is independent from all other variables at depth smaller or equal to the depth of $s$.

Note that, from (CI3), definition 13.42 implies that, for all $s \in G_d$,

$$(X^{(s)} \perp\!\!\!\perp X^{(G_d \setminus \{s\})} \mid X^{(G_q)}, q < d),$$

which, using proposition 13.6, implies that the variables $(X^{(s)}, s \in G_d)$ are mutually independent given $X^{(G_q)}, q < d$. This implies that, for $d = 1$ (letting $s_0$ denote the root in $G$):

$$\mathbb{P}(X^{(G_1)} = x^{(G_1)}, X^{(s_0)} = x^{(s_0)}) = \mathbb{P}(X^{(s_0)} = x^{(s_0)}) \prod_{s \in G_1} \mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(s_0)} = x^{(s_0)}).$$

(If $\mathbb{P}(X^{(s_0)} = x^{(s_0)}) = 0$, the choice for the conditional probabilities can be made arbitrarily without changing the left-hand side which vanishes.) More generally, we have, letting $G_{<d} = G_0 \cup \cdots \cup G_{d-1}$,

$$\mathbb{P}(X^{(G_{\leq d})} = x^{(G_{\leq d})}) = \prod_{s \in G_d} \mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(pa(s))} = x^{(pa(s))}) \mathbb{P}(X^{(G_{<d})} = x^{(G_{<d})})$$

(with again an arbitrary choice for conditional probabilities that are not defined) so that, we obtain, by induction, for $x \in \mathcal{F}(V)$

$$\mathbb{P}(X = x) = \mathbb{P}(X^{(s_0)} = x^{(s_0)}) \prod_{s \neq s_0} p_s(x^{(pa(s))}, x^{(s)}) \tag{13.30}$$

where $p_s(x^{(pa(s))}, x^{(s)}) \triangleq \mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(pa(s))} = x^{(pa(s))})$ are the tree transition probabilities between a parent and a child. So we have the following proposition.

**Proposition 13.43** *A process X is Markov relative to a tree $G = (V, E)$ if and only if there exists a probability distribution $p_0$ on $F_{s_0}$ and a family $(p_{st}, (s, t) \in E)$ such that $p_{st}$ is a transition probability from $F_s$ to $F_t$ and*

$$P_X(x) = p_0(x_{s_0}) \prod_{(s,t) \in E} p_{st}(x^{(s)}, x^{(t)}), \ x \in \mathcal{F}(V). \tag{13.31}$$

We only have proved the "only if" part, but the "if" part is obvious from (13.31). Another property that becomes obvious with this expression is the first part of the following proposition.

**Proposition 13.44** *If a process X is Markov relative to a tree $G = (V, E)$ then it is $G^\flat$ Markov. Conversely, if $G = (V, E)$ is an undirected acyclic graph and X is G-Markov, then X is Markov relative to any tree $\tilde{G}$ such that $\tilde{G}^\flat = G$.*

PROOF To prove the converse part, assume that $G = (V, E)$ is undirected acyclic and that $X$ is $G$-Markov. Take $\tilde{G}$ such that $\tilde{G}^\flat = G$. For $s \in V$ and its parent $pa(s)$ in $\tilde{G}$, the sets $\{s\}$ and $\tilde{G}_{\leq d} \setminus \{s, pa(s)\}$ are separated by $pa(s)$ in $G$. To see this, assume that there exists a $t \in \tilde{G}_{\leq d} \setminus \{s, pa(s)\}$ with a path from $t$ to $s$ that does not pass through $pa(s)$. Then we can complete this path with the path from $t$ to the first common ancestor (in $\tilde{G}$) of $t$ and $s$ and back to $s$ to create a path from $s$ to $s$ that passes only once through $\{pa(s), s\}$ and therefore contains a loop by proposition 13.38.

The $G$-Markov property now implies

$$(X^{(s)} \perp\!\!\!\perp (X^{(\tilde{G}_d \setminus \{s\})}, X^{(\tilde{G}_q \setminus \{pa(s)\})}, q < d) \mid X^{(pa(s))})$$

which proves that $X$ is $\tilde{G}$-Markov. ∎

**Remark 13.45** We see that there is no real gain in generality with passing from undirected to directed graphs when working with trees. This is an important remark, because directionality in graphs is often interpreted as causality. For example, there is a natural causal order in the statements

$$(\text{it rains}) \to (\text{car windshields get wet}) \to (\text{car wipers are on})$$

in the sense that each event can be seen as a logical precursor to the next one. However, because one can pass from this directed chain to an equivalent undirected chain and then back to a equivalent directed tree by choosing any of the three variables as roots, there is no way to infer, from the observation of the joint distribution of the three events (it rains, car windshields get wet, wipers are on), any causal relationship between them: the joint distribution cannot resolve whether wipers are on because

it rains, or whether turning wipers on automatically wets windshields which in turn triggers a shower !

To infer causal relationships, one needs a different kind of observation, that would modify the distribution of the system. Such an operation (called an intervention), can be done, for example, by preventing the windshields from being wet (doing, for example, the observation in a parking garage), or forcing them to be wet (using a hose). Then, one can compare observations made with these new conditions, and those made with the original system, and check, for example, whether they modified the probability that rain occurs outside. The answer (likely to be negative !) would refute any causal relationship from "windshields are wet" to "it rains." On the other hand, the intervention might modify how wipers are used, which would indicate a possible causal relationship from "windshields are wet" to "wipers are on." ♦

## 13.5   Examples of general "loopy" Markov random fields

We will see that acyclic models have very nice computational properties that make them attractive in designing distributions. However, the absence of loops is a very restrictive constraint, which is not realistic in many practical situations. Feedback effects are often needed, for example. Most models in statistical physics are supported by a lattice, in which natural translation/rotation invariance relations forbid using any non-trivial acyclic model. As an example, we now consider the 2D Ising model on a finite grid, which is a model for (anti)-ferromagnetic interaction in a spin system.

Let $G = (V, E)$. A (positive) $G$-Markov model is said to have only pair interactions if and only if can be written in the form

$$\pi(x) = \frac{1}{Z} \exp\left( -\sum_{s \in G} h_s(x^{(s)}) - \sum_{\{s,t\} \in E} h_{\{s,t\}}(x^{(s,t)}) \right).$$

Relating to theorem 13.30, this says that $\pi$ is associated to a potential involving cliques of order 2 at most (note that this does not mean that the cliques of the associated graph have order 2 at most; there can be higher-order cliques, which would then have a zero potential). The functions in the potential are indexed by sets, as they should be from the general definition. However, models with pair interactions are often written in the form

$$\pi(x) = \frac{1}{Z} \exp\left( -\sum_{s \in G} h_s(x^{(s)}) - \sum_{\{s,t\} \in E} \tilde{h}_{st}(x^{(s)}, x^{(t)}) \right)$$

with $\tilde{h}_{st}(\lambda, \mu) = \tilde{h}_{ts}(\mu, \lambda)$ (which is equivalent, taking $\tilde{h} = h/2$).

The *Ising model* is a special case of models with pair interactions, for which the state space, $F_s$, is equal to $\{-1, 1\}$ for all $s$ and

$$h_s(x^{(s)}) = \alpha_s x^{(s)}, \quad h_{\{s,t\}}(x^{(s)}, x^{(t)}) = \beta_{st} x^{(s)} x^{(t)}.$$

In fact, for binary variables, this is the most general pair interaction model.



Figure 13.3:  Graph forming a two-dimensional regular grid.

The Ising model is moreover usually defined on a regular lattice, which, in two dimensions, implies that $V$ is a finite rectangle in $\mathbb{Z}^2$, for example $V = \{-N, \ldots, N\}^2$. The simplest choice of a translation- and 90-degree rotation-invariant graph is the nearest-neighbor graph for which $\{(i, j), (i', j')\} \in E$ if and only if $|i - i'| + |j - j'| = 1$ (see fig. 13.3). With this graph, one can furthermore simplify the model to obtain the *isotropic Ising model* given by

$$\pi(x) = \frac{1}{Z} \exp\left(-\alpha \sum_{s \in V} x^{(s)} - \beta \sum_{s \sim t} x^{(s)} x^{(t)}\right).$$

When $\beta < 0$, the model is *ferromagnetic*: each pair of neighbors with identical signs brings a negative contribution to the energy, making the configuration more likely (since lower energy implies higher probability).

The Potts model generalizes the Ising model to finite, but non-necessarily binary, state spaces, say, $F_s = F = \{1,\ldots,n\}$. Define the function $\delta(\lambda,\mu) = 1$ if $\lambda = \mu$ and $(-1)$ otherwise. Then the Potts model is given by

$$\pi(x) = \frac{1}{Z}\exp\left(-\alpha\sum_{s\in V}h(x^{(s)}) - \beta\sum_{s\sim t}\delta(x^{(s)},x^{(t)})\right) \tag{13.32}$$

for some function $h$ defined on $F$.

## 13.6   General state spaces

Our discussion of Markov random fields on graphs was done under the assumption of finite state spaces, which notably simplifies many of the arguments and avoids relying too much on measure theory. While this situation does cover a large range of application, there are cases in which one wants to consider variables taking values in continuous spaces, or in countable (infinite) spaces.

The results obtained for discrete variables can most of the time be extended to variables whose distribution has a p.d.f. with respect to a product of measures on the sets in which they take their values. For example, let $X, Y, Z$ takes values in $\mathcal{R}_X, \mathcal{R}_Y, \mathcal{R}_Z$, equipped with $\sigma$-algebras $\mathcal{S}_X$, $\mathcal{S}_Y$, $\mathcal{S}_Z$ and measures $\mu_X$, $\mu_Y$, $\mu_Z$. Assume that $P_{X,Y,Z}$ is absolutely continuous with respect to $\mu_X \otimes \mu_Y \otimes \mu_Z$, with density $\varphi_{XYZ}$. In such a situation, (13.3) remains valid, in that $X$ is conditionally independent of $Y$ given $Z$ if and only if

$$\varphi_{XYZ}(x,y,z)\varphi_Z(z) = \varphi_{XZ}(x,z)\varphi_{YZ}(y,z) \tag{13.33}$$

almost everywhere (relative to $\mu_X \otimes \mu_Y \otimes \mu_Z$). Here, $\varphi_{XZ}, \varphi_{YZ}, \varphi_Z$ are marginal densities of the indexed random variables. The only difficulty in the argument, provided below for the interested reader, is dealing properly with sets of measure zero.

PROOF (PROOF OF (13.33))  Introduce the conditional densities

$$\varphi_{XY|Z}(x,y \mid z) = \frac{\varphi_{XYZ}(x,y,z)}{\varphi_Z(z)}$$

and similarly $\varphi_{X|Z}$ and $\varphi_{Y|Z}$, which are defined when $z \notin M_Z = \{z \in \mathcal{R}_Z : \varphi_Z(z) = 0\}$. By definition of conditional independence, we have, for all $A \in \mathcal{S}_X$, $B \in \mathcal{S}_X$

$$\int_{A\times B}\varphi_{XY|Z}(x,y \mid z)\mu_X(dx)\mu_Y(dy) = \int_{A\times B}\varphi_{X|Z}(x \mid z)\varphi_{Y|Z}(y \mid z)\mu_X(dx)\mu_Y(dy)$$

for all $z \notin M_Z$, which implies that, for all $z \notin M_Z$, there exists a set $N_z \subset R_X \times R_Y$ such that $\mu_X \times \mu_Y(N_z) = 0$ and

$$\varphi_{XY|Z}(x,y \mid z) = \varphi_{X|Z}(x \mid z)\varphi_{Y|Z}(y \mid z)$$

for all $z \notin M_Z$ and $(x,y) \notin N_z$. This immediately implies (13.33) for those $(x,y,z)$.

If $z \in M_Z$, then

$$0 = \varphi_Z(z) = \int_{\mathcal{R}_X} \varphi_{XZ}(x,z)\mu_X(dx) = \int_{\mathcal{R}_Y} \varphi_{YZ}(x,z)\mu_Y(dy)$$

implying that $\varphi_{XZ}(x,z) = \varphi_{YZ}(y,z) = 0$ excepted on some set $N_z$ such that $\mu_X \otimes \mu_Y(N_z) = 0$, and (13.33) is therefore also true outside of this set. Now, letting $N = \{(x,y,z) : (x,y) \in N_z\}$, we find that (13.33) is true for all $(x,y,z) \notin N$ and

$$\mu_X \otimes \mu_Y \otimes \mu_Z(N) = \int_{\mathcal{R}_X \times \mathcal{R}_Y \times \mathcal{R}_Z} \mathbf{1}_{(x,y) \in N_z} \mu_X(dx)\mu_Y(dy)\mu_Z(dz) = \int_{\mathcal{R}_Z} \mu_X \otimes \mu_Y(N_z)\mu_Z(dz) = 0.$$

(This argument involves Fubini's theorem [171].)  ■

With this definition, the proof of proposition 13.5 can be caried on without change, with the positivity condition expressing the fact that there exists $\tilde{R}_X \subset \mathcal{R}_X$, $\tilde{R}_Y \subset \mathcal{R}_Y$ and $\tilde{R}_Z \subset \mathcal{R}_X$ such that $\varphi_{XYZ}(x,y,z) > 0$ for all $x,y,z \in \tilde{R}_X \times \tilde{R}_Y \times \tilde{R}_Z$. (This proposition is actually valid in full generality, with a proper definition of positivity.)

When considering random fields with general state spaces, we will restrict to the similar situation in which each state space $F_s$ is equipped with a $\sigma$-algebra $\mathcal{S}_s$ and a measure $\mu_s$, and the joint distribution, $P_X$ of the random field $X = (X_s, s \in V)$ is absolutely continuous with respect to $\mu \stackrel{\Delta}{=} \bigotimes_{s \in V} \mu_s$, denoting by $\pi$ the corresponding p.d.f. We will says that $\pi$ is positive if there exists $\tilde{F} = (\tilde{F}_s, s \in V)$ with measurable $\tilde{F}_s \subset F_s$ such that $\pi(x) > 0$ for all $x \in \mathcal{F}(V, \tilde{F})$. Without loss of generality unless one considers multiple random fields with different supports, we will assume that $\tilde{F}_s = F_s$ for all $s$.

The definition of consistent families of local interactions (definition 13.24) must be modified by adding the condition that

$$\int_{\mathcal{F}(V)} \prod_{C \in \mathcal{C}} \varphi_C(x^{(C)})\mu(dx) < \infty. \tag{13.34}$$

This requirement is obviously needed to ensure that the normalizing constant in (13.21) is finite. Proposition 13.25 is then true (with sums replaced by integrals in the proof) and so are propositions 13.26 and 13.27. Finally, the Hammersley-Clifford theorem (theorem 13.29) extends to this context.

Even though it is a natural requirement, condition (13.34) may be hard to assess with general families of local interactions. In the case of Gaussian distributions, however, one can provide relatively simple conditions. Assume that $F_s = \mathbb{R}$ for all

$s \in V$, and condider a potential $\Lambda = (\lambda_C, C \in \mathcal{C})$ with only univariate and bivariate interactions, such that, for some vector $a \in \mathbb{R}^d$ (with $d = |V|$) and symmetric matrix $b \in \mathcal{S}_d$,

$$\begin{cases} \lambda_{\{s\}}(x^{(\{s\})}) = -a^{(s)}x^{(s)} + \dfrac{1}{2}b_{ss}(x^{(s)})^2 \\ \lambda_{\{s,t\}}(x^{(\{s,t\})}) = b_{st}x^{(s)}x^{(t)} \end{cases}$$

Then, considering $x \in \mathcal{F}(V)$ as a $d$-dimensional vector, we have

$$\pi(x) = \frac{1}{Z}\exp\left(a^T x - \frac{1}{2}x^T b x\right),$$

with the integrability requirement that $b > 0$ (positive definite). The random field then follows a Gaussian distribution with mean $m = b^{-1}a$ and covariance matrix $\Sigma = b^{-1}$. The normalizing constant, $Z$, is given by

$$Z = \frac{e^{-\frac{1}{2}a^T b a}(2\pi)^{d/2}}{\sqrt{\det b}}.$$

This Markov random field parametrization of Gaussian distributions emphasizes the conditional structure of the variables rather than their covariances. It is useful when the associated graph, represented by the matrix $b$ is sparse. In particular, the conditional distribution of $X^{(s)}$ given the other variables is Gaussian, with mean $(a^{(s)} - \sum_{t \neq s} b_{st}x^{(t)})/b_{ss}$ and variance $1/b_{ss}$.

# Chapter 14

# Probabilistic Inference for Random Fields

Once the joint distribution of a family of variables has been modeled as a random field, this model can be used to estimate the probabilities of specific events, or the expectations of random variables of interest. For example, if the modeled variables relate to a medical condition, in which variables such as diagnosis, age, gender, clinical evidence can interact, one may want to compute, say, the probability of someone having a disease given other observable factors. Note that, being able to compute expectations of the modeled variables for $G$-Markov processes also ensures that one can compute conditional expectations of some modeled variables given others, since, by proposition 13.22, conditional $G$-Markov distributions are Markov over restricted graphs.

We assume that $X$ is $G$-Markov for a graph $G = (V, E)$ and restrict (unless specified otherwise) to finite state spaces. We condider the basic problem to compute $\mathbb{P}(X^{(S)} = x^{(S)})$ when $S \subset V$, starting with one-vertex marginals, $\mathbb{P}(X^{(s)} = x^{(s)})$.

The Hammersley-Clifford theorem provides a generic form for general positive $G$-Markov processes, in the form

$$\mathbb{P}(X = x) = \pi(x) = \frac{1}{Z} \exp\left(-\sum_{C \in \mathcal{C}_G} h_C(x^{(C)})\right). \tag{14.1}$$

So, formally, marginal distributions are given by the ratio

$$\mathbb{P}(X^{(S)} = x^{(S)}) = \frac{\sum_{y \in \mathcal{F}(V), y^{(S)} = x^{(S)}} \exp\left(-\sum_{C \in \mathcal{C}_G} h_C(y^{(C)})\right)}{\sum_{y \in \mathcal{F}(V)} \exp\left(-\sum_{C \in \mathcal{C}_G} h_C(y^{(C)})\right)}.$$

The problem is that the sums involved in this ratio involve a number of terms that grows exponentially with the size of $V$. Unless $V$ is very small, a direct computation of these sums is intractable. An exception to this is the case of acyclic graphs, as

we will see in section 14.2. But for general, loopy, graphs, the sums can only be approximated, using, for example, Monte-Carlo sampling, as described in the next section.

## 14.1   Monte Carlo sampling

Markov chain Monte Carlo methods are well adapted to sampling from Markov random fields, because conditional distributions used in Gibbs sampling, or, more generally, ratios of probabilities used in the Metropolis-Hastings algorithm do not in require the computation of the normalizing constant $Z$ in (14.1). The simplest use of Gibbs sampling generalizes the Ising model example of section 12.4.2. Using the notation of Algorithm 12.2, one lets $\mathcal{B}'_s = \mathcal{F}(s^c)$ (with the notation $s^c = V \setminus \{s\}$) and $U_s(x) = x^{(s^c)}$. The conditional distribution given $U_s$ is

$$Q_s(U_s(x), y) = \mathbb{P}(X^{(s)} = y^{(s)} \mid X^{(s^c)} = x^{(s^c)})\mathbf{1}_{y^{(s^c)}=x^{(s^c)}}.$$

The conditional probability in the r.h.s. of this equation takes the form

$$\pi_s(y^{(s)} \mid x^{(s^c)}) \stackrel{\Delta}{=} \mathbb{P}(X^{(s)} = y^{(s)} \mid X^{(s^c)} = x^{(s^c)}) = \frac{1}{Z_s(x^{(s^c)})} \exp\left(- \sum_{C \in \mathcal{C}, s \in C} h_C(y^{(s)} \wedge x^{(C \cap s^c)})\right)$$

with

$$Z_s(x^{(s^c)}) = \sum_{z^{(s)} \in F_s} \exp\left(- \sum_{C \in \mathcal{C}, s \in C} h_C(z^{(s)} \wedge x^{(C \cap s^c)})\right).$$

The Gibbs sampling algorithm samples from $Q_s$ by visiting all $s \in V$ infinitely often, as described in Algorithm 12.2. Metropolis-Hastings schemes are implemented similarly, the most common choice using a local update scheme in Algorithm 12.3 such that $g(x, \cdot)$ only changes one coordinate, chosen at random, so that

$$g(x, y) = \frac{1}{|V|} \sum_{s \in V} \mathbf{1}_{y^{(s^c)}=x^{(s^c)}} g_s(y^{(s)})$$

where $g_s$ is some probability distribution on $F_s$. The acceptance probability $a(x, y)$ is equal to 1 when $y = x$. If $y \neq x$ and $g(x, y) > 0$, there is a unique $s$ for which $y^{(s^c)} = x^{(s^c)}$ and

$$a(x, y) = \min\left(1, \frac{\pi(y)g(y, x)}{\pi(x)g(x, y)}\right)$$

with

$$\frac{\pi(y)g(y, x)}{\pi(x)g(x, y)} = \frac{\pi_s(y^{(s)} \mid x^{(s^c)})g_s(x^{(s)})}{\pi_s(x^{(s)} \mid x^{(s^c)})g_s(y^{(s)})}.$$

Note that the latter equation avoids the computation of the local normalizing constant $Z_s(x^{(s^c)})$, which simplifies in the ratio.

Both algorithms have a transition probability $P$ that satisfies $P^m(x, y) > 0$ for all $x, y \in \mathcal{F}(V)$, with $m = |V|$ (for Metropolis-Hastings, one must assume that $g_s(y^{(s)}) > 0$ for all $y^{(s)} \in F_s$. This ensures that the chain is uniformly geometrically ergodic, i.e., (12.10) is satisfied with a constant $M$ and some $\rho < 1$. However, in many practical cases (especially for strongly structured distributions and large sets $V$), the convergence rate, $\rho$ can be very close to 1, resulting in a slow convergence.

Acceleration strategies have been designed to address this issue, which is often due to the existence of multiple configurations that are local modes of the probability $\pi$. Such configurations are isolated from other high-probability configurations because local updating schemes need to make multiple low-probability changes to access them from the local mode. The following two approaches provide examples designed to address this issue.

**a. Cluster sampling.** To facilitate escaping from such local modes, it is sometimes possible to augment the state space by introducing a new configuration space, with variable denoted $\xi$, and designing a joint distributions $\hat{\pi}(\xi, x)$ such that the marginal distribution on $\mathcal{F}(V)$ (summing over $\xi$) is the targeted $\pi$. The additional variable can create high-probability bridges between local modes for $\pi$, and accelerate convergence.

To take an example, assume that all sets $F_s$ are identical (letting $F = F_s$, $s \in V$) and that the auxiliary variable $\xi$ takes values in the set of functions from $E$ to $\{0, 1\}$, that we will denote $\mathcal{B}(E)$, i.e., that it takes the form $(\xi^{(st)}, \{s, t\} \in E)$, with $\xi^{(st)} \in \{0, 1\}$. For $x \in \mathcal{F}(V)$, introduce the set $\mathcal{B}_x$ containing all $\xi \in \mathcal{B}(E)$, such that for all $\{s, t\} \in E$,

$$x^{(s)} \neq x^{(t)} \Longrightarrow \xi^{(st)} = 1.$$

Assume that the conditional distribution of $\xi$ given $x$ is supported by $\mathcal{B}_x$, such that, for $\xi \in \mathcal{B}_x$

$$\mathbb{P}(\xi = \xi \mid X = x) = \hat{\pi}(\xi \mid x) = \frac{1}{\zeta(x)} \exp\left(-\sum_{\{s,t\} \in E} \mu_{st} \xi^{(st)}\right).$$

The coefficients $\mu_{st}$ are free to choose (and one possible choice is to take $\mu_{st} = 0$ for all $\{s, t\} \in E$). For this distribution, all $\xi^{(st)}$ are independent conditionally to $X = x$, with $\xi^{(st)} = 1$ with probability 1 if $x^{(s)} \neq x^{(t)}$, and

$$P(\xi^{(st)} = 1 \mid X = x) = \frac{e^{-\mu_{st}}}{1 + e^{-\mu_{st}}} \tag{14.2}$$

if $x^{(s)} = x^{(t)}$. This conditional distribution is, as a consequence, very easy to sample

from. Moreover, the normalizing constant $\zeta(x)$ has closed form and is given by

$$\zeta(x) = \prod_{\{s,t\}\in E} (\mathbf{1}_{x^{(s)}=x^{(t)}} + e^{-\mu_{st}}) = \exp\left(\sum_{\{s,t\}\in E} \log(1 + e^{-\mu_{st}}) + \sum_{\{s,t\}\in E} \log(1 + e^{\mu_{st}})\mathbf{1}_{x^{(s)}\neq x^{(t)}}\right).$$

Now consider the conditional probability that $X = x$ given $\boldsymbol{\xi} = \xi$. For this distribution, one has, with probability 1, $X^{(s)} = X^{(t)}$ when $\xi^{(st)} = 0$. This implies that $X$ is constant on the connected components of the subgraph $(V, E_\xi)$ of $(V, E)$, where $\{s, t\} \in E_\xi$ if and only if $\xi^{(st)} = 0$. Let $V_1, \ldots, V_m$ denote these connected components (these components and their number depend on $\xi$). The conditional distribution of $X$ given $\xi$ is therefore supported by the configurations such that there exists $c_1, \ldots, c_m \in F$ such that $x^{(s)} = c_j$ if and only if $s \in V_j$, that we will denote, with some abuse of notation: $c_1^{(V_1)} \wedge \cdots \wedge c_m^{(V_m)}$.

Given this remark, the conditional distribution of $X$ given $\boldsymbol{\xi} = \xi$ is equivalent to a distribution on $F^m$, which may be feasible to sample from directly if $|F|$ and $m$ are not too large. To sample from $\pi$, one now needs to alternate between sampling $\xi$ given $X$ and the converse, yielding the following first version of cluster-based sampling.

---

**Algorithm 14.1 (Cluster-based sampling: Version 1)**
This algorithm samples from (14.1).

(1) Initialize the algorithm some configuration $x \in \mathcal{F}(V)$.

(2) Loop over the following steps:

    a.  Generate a configuration $\xi \in \mathcal{B}_x$ such that $\xi^{(st)} = 1$ with probability given by (14.2) when $x^{(s)} = x^{(t)}$.

    b.  Determine the connected components, $V_1, \ldots, V_m$, of the graph $G_\xi = (V, E_\xi)$ with edges given by pairs $\{s, t\}$ such that $\xi^{(st)} = 1$.

    c.  Sample values $c_1, \ldots, c_m \in F$ according to the distribution

$$q(c_1, \ldots, c_m) \propto \frac{\pi(c_1^{(V_1)} \wedge \cdots \wedge c_m^{(V_m)})}{\zeta(c_1^{(V_1)} \wedge \cdots \wedge c_m^{(V_m)})}.$$

    d.  Set $x = c_1^{(V_1)} \wedge \cdots \wedge c_m^{(V_m)}$.

---

Step (2.c) takes a simple form in the special case when $\pi$ is a non-homogeneous Potts model ((13.32)) with positive interactions, that we will write as

$$\pi(x) = \exp\left(-\sum_{s\in V} \alpha_s x^{(s)} - \sum_{\{s,t\}\in E} \beta_{st}\mathbf{1}_{x^{(s)}\neq x^{(t)}}\right)$$

with $\beta_{st} \geq 0$. Then

$$\frac{\pi(x)}{\zeta(x)} \propto \exp\left(-\sum_{s \in V} \alpha_s x^{(s)} - \sum_{\{s,t\} \in E} (\beta_{st} - \beta'_{st})\mathbf{1}_{x_s \neq s_t}\right)$$

with $\beta'_{st} = \log(1 + e^{\mu_{st}})$. If one chooses $\mu_{st}$ such that $\beta'_{st} = \beta_{st}$ (which is possible since $\beta_{st} \geq 0$), then the interaction term disappears and the probability $q$ in (2.c) is proportional to

$$\prod_{j=1}^{m} \exp\left(-\sum_{s \in V_j} \alpha_s\right)$$

so that $c_1, \ldots, c_m$ can be generated independently. The resulting algorithm is the Swendsen-Wang sampling algorithm for the Potts model [186]. The presentation given here adapts the one introduced in Barbu and Zhu [16].

For more general models, step (2.c) can be computationally costly, especially if the number of connected components is large. In this case, this step can be replaced by a Gibbs sampling step for one of the $c'_j s$ conditional to the others (and $\xi$) that we summarize in the following variation of Algorithm 14.1.

---

**Algorithm 14.2 (Cluster-based sampling: Version 2)**
This algorithm samples from (14.1).

(1) Initialize the algorithm some configuration $x \in \mathcal{F}(V)$.

(2) Loop over the following steps:

    a. Generate a configuration $\xi \in \mathcal{B}_x$ such that $\xi^{(st)} = 1$ with probability given by (14.2) when $x^{(s)} = x^{(t)}$.

    b. Determine the connected components, $V_1, \ldots, V_m$, of the graph $G_\xi = (V, E_\xi)$ with edges given by pairs $\{s, t\}$ such that $\xi^{(st)} = 1$. Note that $x$ is constant on each of these connected components, i.e., there exists $c_1, \ldots, c_m \in F$ such that $x = c_1^{(V_1)} \wedge \cdots \wedge c_m^{(V_m)}$.

    c. Select at random one of the components, say, $j_0 \in \{1, \ldots, m\}$.

    d. Sample the value $\tilde{c}_{j_0} \in F$ according to the distribution

$$q(\tilde{c}_{j_0}) \propto \frac{\pi(\tilde{c}_1^{(V_1)} \wedge \cdots \wedge \tilde{c}_m^{(V_m)})}{\zeta(\tilde{c}_1^{(V_1)} \wedge \cdots \wedge \tilde{c}_m^{(V_m)})}.$$

       with $\tilde{c}_j = c_j$ if $j \neq j_0$.

    e. Set $x^{(s)} = \tilde{c}_{j_0}$ for $s \in V_{j_0}$.

---

Unlike single-variable updating schemes, these algorithms can update large chunks of the configurations at each step, and may result in significantly faster convergence

of the sampling procedure. Note that step (2.d) in Algorithm 14.2 can be replaced by a Metropolis-Hastings update with a proper choice of proposal probability [16].

**b. Parallel tempering.** We now consider a different kind of extension in which we allow $\pi$ depends continuously on a parameter $\beta > 0$, writing $\pi_\beta$ and, the goal is to sample from $\pi_1$. For example, one can extend (14.1) by the family of probability distributions

$$\pi_\beta(x) = \frac{1}{Z_\beta} \exp\left(-\beta \sum_{C \in \mathcal{C}_G} h_C(x^{(C)})\right)$$

for $\beta \geq 0$. For small $\beta$, $\pi_\beta$ gets close to the uniform distribution on $\mathcal{F}(V)$ (achieved for $\beta = 0$), so that it becomes easier to move from local mode to local mode. This implies that sampling with small $\beta$ is more efficient and the associated Markov chain moves more rapidly in the configuration space.

Assume given, for all $\beta$, two ergodic transition probabilities on $\mathcal{F}(V)$, $q_\beta$ and $\tilde{q}_\beta$ such that (12.7) is satisfied with $\pi_\beta$ as invariant probability, namely

$$\pi_\beta(y)q_\beta(y,x) = \pi_\beta(x)\tilde{q}_\beta(x,y) \tag{14.3}$$

for all $x, y \in \mathcal{F}(V)$ (as seen in (12.7), $\tilde{q}_\beta$ is the transition probability for the reversed chain). The basic idea is that $q_\beta$ provides a Markov chain that converges rapidly for small $\beta$ and slowly when $\beta$ is closer to 1. Parallel tempering (this algorithm was introduced in Neal [140] based on ideas developed in Marinari and Parisi [126]) leverages this fact (and the continuity of $\pi_\beta$ in $\beta$) to accelerate the simulation of $\pi_1$ by introducing intermediate steps sampling at low $\beta$ values.

The algorithm specifies a sequence of parameters $0 \leq \beta_1 \leq \cdots \leq \beta_m = 1$. One simulation steps goes down, then up this scale, as described in the following algorithm.

---

**Algorithm 14.3 (Parallel Tempering)**
Start with an initial configuration $x_0 \in \mathcal{F}(V)$. This configuration is then updated at each step, using the following sequence of operations.

(1) For $j = 1, \ldots, m$, generate a configuration $x_j$ according to $\tilde{q}_{\beta_j}(x_{j-1}, \cdot)$.

(2) Generate a configuration $z_{m-1}$ according to $q_{\beta_m}(x_m, \cdot)$.

(3) For $j = m - 1, \ldots, 1$, generate a configuration $z_{j-1}$ according to $q_{\beta_j}(z_j, \cdot)$.

(4) Set $x_0 = z_0$ with probability

$$\min\left(1, \frac{\pi_{\beta_0}(z_0)}{\pi_{\beta_0}(x_0)}\left(\prod_{j=1}^{m-1} \frac{\pi_{\beta_j}(x_{j-1})}{\pi_{\beta_j}(x_j)}\right)\frac{\pi_{\beta_m}(x_{m-1})}{\pi_{\beta_m}(z_{m-1})}\left(\prod_{j=1}^{m-1} \frac{\pi_{\beta_j}(z_j)}{\pi_{\beta_j}(z_{j-1})}\right)\right).$$

(Otherwise, keep $x_0$ unchanged).

Importantly, the acceptance probability at step (4) only involves ratios of $\pi'_\beta s$ and therefore no normalizing constant. We now show that this algorithm is $\pi_{\beta_0}$-reversible. Let $p(\cdot,\cdot)$ denote the transition probability of the chain. If $z_0 \neq x_0$, $p(x_0, z_0)$ corresponds to steps (1) to (3), with acceptance at step(4), and is therefore given by the sum, over all $x_1, \ldots, x_m$ and $z_1, \ldots, z_m$, of products

$$\tilde{q}_{\beta_1}(x_0, x_1) \cdots \tilde{q}_{\beta_m}(x_{m-1}, x_m) q_{\beta_m}(x_m, z_{m-1}) \cdots q_{\beta_1}(z_1, z_0)$$

$$\min\left(1, \frac{\pi_{\beta_0}(z_0)}{\pi_{\beta_0}(x_0)} \left(\prod_{j=1}^{m-1} \frac{\pi_{\beta_j}(x_{j-1})}{\pi_{\beta_j}(x_j)}\right) \frac{\pi_{\beta_m}(x_{m-1})}{\pi_{\beta_m}(z_{m-1})} \left(\prod_{j=1}^{m-1} \frac{\pi_{\beta_j}(z_j)}{\pi_{\beta_j}(z_{j-1})}\right)\right)$$

Applying (14.3), this is equal to

$$\min\Big(\tilde{q}_{\beta_1}(x_0, x_1) \cdots \tilde{q}_{\beta_m}(x_{m-1}, x_m) q_{\beta_m}(x_m, z_{m-1}) \cdots q_{\beta_1}(z_1, z_0),$$

$$\frac{\pi_{\beta_0}(z_0)}{\pi_{\beta_0}(x_0)} q_{\beta_1}(x_0, x_1) \cdots q_{\beta_m}(x_{m-1}, x_m) \tilde{q}_{\beta_m}(x_m, z_{n-1}) \cdots \tilde{q}_{\beta_1}(z_1, z_0)\Big)$$

So,

$$\pi_{\beta_0}(x_0) p(x_0, z_0) = \sum \min\Big(\pi_{\beta_0}(x_0) \tilde{q}_{\beta_1}(x_0, x_1) \cdots \tilde{q}_{\beta_m}(x_{m-1}, x_m) q_{\beta_m}(x_m, z_{m-1}) \cdots q_{\beta_1}(z_1, z_0),$$

$$\pi_{\beta_0}(z_0) q_{\beta_1}(x_1, x_0) \cdots q_{\beta_m}(x_m, x_{m-1}) \tilde{q}_{\beta_m}(z_{m-1}, x_m) \cdots \tilde{q}_{\beta_1}(z_0, z_1)\Big)$$

where the sum is over all $x_1, \ldots, x_m, z_1, \ldots, z_{m-1} \in \mathcal{F}(V)$. The sum is, of course, unchanged if one renames $x_1, \ldots, x_m, z_1, \ldots, z_{m-1}$ to $z_1, \ldots, z_m, x_1, \ldots, x_{m-1}$, but doing so provides the expression of $\pi_{\beta_0}(z_0) p(z_0, x_0)$, proving the reversibility of the chain with respect to $\pi_{\beta_0}$.

## 14.2  Inference with acyclic graphs

We now switch to deterministic methods to compute, or approximate, marginal probabilities of Markov random fields. In this section, we consider a directed acyclic graph $G = (V, E)$. As we have seen, Markov processes for acyclic graphs are also Markov for any tree structure associated with the graph. Introducing such a tree, $\tilde{G} = (V, \tilde{E})$ with $\tilde{G}^\flat = G$, we know that a Markov process on $G$ can be written in the form (letting $s_0$ denote the root in $\tilde{G}$):

$$\pi(x) = p_{s_0}(x^{(s_0)}) \prod_{(s,t) \in \tilde{E}} p_{st}(x^{(s)}, x^{(t)}) \tag{14.4}$$

where $p_{s_0}$ is a probability and $p_{st}$ a transition probability.

We now show how to compute marginal probabilities of configurations $x^{(S)}$, denoted $\pi_S(x^{(S)})$, for a set $S \subset V$, starting with singletons $S = \{s\}$. The computation can be done by propagating down the tree as follows. For $s = s_0$, the probability is known, with $\pi_{s_0} = p_{s_0}$. Now take an arbitrary $s \neq s_0$ and let $pa(s)$ be its parent. Then

$$\pi_s(x^{(s)}) = \mathbb{P}(X^{(s)} = x^{(s)}) = \sum_{y^{(pa(s))} \in F_{pa(s)}} P(X^{(s)} = x^{(s)} \mid X^{(pa(s))} = y^{(pa(s))}) P(x^{(pa(s))} = y^{(pa(s))})$$

$$= \sum_{y^{(pa(s))} \in F_{pa(s)}} \pi_{pa(s)}(y^{(pa(s))}) p_{pa(s)}(y_{pa(s)}, x^{(s)})$$

so that the marginal probability at any $s \neq s_0$ can be computed given the marginal probability of its parent. We can propagate the computation down the tree, with a total cost for computing $\pi_s$ proportional to $\sum_{k=1}^{n} |F_{t_{k-1}}||F_{t_k}|$ where $t_0 = s_0, t_1, \ldots, t_n = s$ is the unique path between $s_0$ and $s$. This is linear in the depth of the tree, and quadratic (not exponential) in the sizes of the state spaces. The computation of all singleton marginals requires an order of $\sum_{(s,t) \in E} |F_s||F_t|$ operations.

Now, assume that probabilities of singletons have been computed and consider an arbitrary set $S \subset V$. Let $s \in V$ be an ancestor of every vertex in $S$, maximal in the sense that none of its children also satisfy this property. Consider the subtrees of $\tilde{G}$ starting from each of the children of $s$, denoted $\tilde{G}_1, \ldots, \tilde{G}_n$ with $\tilde{G}_k = (V_k, \tilde{E}_k)$. Let $S_k = S \cap V_k$. From the conditional independence,

$$\pi_S(x^{(S)}) = \sum_{y^{(s)} \in F_s} P(X^{(S \setminus \{s\})} = x^{(S \setminus \{s\})} \mid X^{(s)} = y^{(s)}) \pi_s(y^{(s)})$$

$$= \sum_{y^{(s)} \in F_s} \prod_{k=1, S_k \neq \emptyset}^{n} P(X^{(S_k)} = x^{(S_k)} \mid X^{(s)} = y^{(s)}) \pi_s(y_s)$$

Now, for all $k = 1, \ldots, n$, we have $|S_k| < |S|$: this is obvious if $S$ is not completely included in one of the $V_k$'s. But if $S \subset V_k$ then the root, $s_k$, of $V_k$ is an ancestor of all the elements in $S$ and is a child of $s$, which contradicts the assumption that $s$ is maximal. So we have reduced the computation of $\pi_S(x_S)$ to the computations of $n$ probabilities of smaller sets, namely $P(X^{(S_k)} = x^{(S_k)} \mid X^{(s)} = y^{(s)})$ for $S_k \neq \emptyset$. Because the distribution of $X^{(V_k)}$ conditioned at $s$ is a $\tilde{G}_k$-Markov model, we can reiterate the procedure until only sets of cardinality one remain, for which we know how to explicitly compute probabilities.

This provides a feasible algorithm to compute marginal probabilities with trees, at least when its distribution is given in tree-form, like in (14.4). We now address

the situation in which one starts with a probability distribution associated with pair interactions (cf. definition 13.24) over the acyclic graph $G$

$$\pi(x) = \frac{1}{Z} \prod_{s \in V} \varphi_s(x^{(s)}) \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}). \tag{14.5}$$

We assume these local interactions to be consistent, still allowing for some vanishing $\varphi_{st}(x^{(s)}, x^{(t)})$.

Putting $\pi$ in the form (14.4) is equivalent to computing all joint probability distributions $\pi_{st}(x^{(s)}, x^{(t)})$ for $\{s, t\} \in E$, and we now describe this computation. Denote

$$U(x) = \prod_{s \in V} \varphi_s(x^{(s)}) \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)})$$

so that $Z = \sum_{y \in \mathcal{F}(V)} U(y)$. For the tree $\tilde{G} = (V, \tilde{E})$, and $t \in V$, we let $\tilde{G}_t = (V_t, \tilde{E}_t)$ be the subtree of $G$ rooted at $t$ (containing $t$ and all its descendants). For $S \subset V$, define

$$U_S(x^{(S)}) = \prod_{s \in S} \varphi_s(x^{(s)}) \prod_{\{s,s'\} \in E, s, s' \in D} \varphi_{ss'}(x^{(s)}, x^{(s')})$$

and

$$Z_t(x^{(t)}) = \sum_{y^{(V_t^*)} \in \mathcal{F}(V_t^*)} U_{V_t}(x^{(t)} \wedge y^{(V_t^*)}).$$

with $V_t^* = V_t \setminus \{t\}$.

**Lemma 14.1** *Let $G = (V, E)$ be a directed acyclic graph and $\pi = P^X$ be the $G$-Markov distribution given by (14.5). With the notation above, we have*

$$\pi_{s_0}(x^{(s_0)}) = \frac{Z_{s_0}(x^{(s_0)})}{\sum_{y^{(s_0)} \in F_{s_0}} Z_{s_0}(y^{(s_0)})} \tag{14.6}$$

*and, for $(s, t) \in \tilde{E}$,*

$$p_{st}(x^{(s)}, x^{(t)}) = P(X^{(t)} = x^{(t)} \mid X^{(s)} = x^{(s)}) = \frac{\varphi_{st}(x^{(s)}, x^{(t)}) Z_t(x^{(t)})}{\sum_{y^{(t)} \in F_t} \varphi_{st}(x^{(s)}, y^{(t)}) Z_t(y^{(t)})} \tag{14.7}$$

PROOF Let $W_t = V \setminus V_t$. Clearly, $Z = \sum_{x^{(0)} \in F_{s_0}} Z_{s_0}(x^{(0)})$ and $\pi_{s_0}(x^{(0)}) = Z_{s_0}(x^{(0)})/Z$ which gives (14.6). Moreover, if $s \in V$, we have

$$\mathbb{P}(X^{(V_s^*)} = x^{(V_s^*)} \mid X^{(s)} = x^{(s)}) = \frac{\sum_{y^{(W_s)}} U(x^{(V_s)} \wedge y^{(W_s)})}{\sum_{y^{(V_s^*)}, y^{(W_s)}} U(x^{(s)} \wedge y^{(V_s^*)} \wedge y^{(W_s)})}.$$

We can write

$$U(x^{(s)} \wedge y^{(V_s^*)} \wedge y^{(W_s)}) = U_{V_s}(x^{(s)} \wedge y^{(V_s^*)})U_{\{s\}\cup W_s}(x^{(s)} \wedge y^{(W_s)})\varphi_s(x^{(s)})^{-1}$$

yielding the simplified expression

$$\mathbb{P}(X^{(V_s^*)} = x^{(V_s^*)} \mid X^{(s)} = x^{(s)}) = \frac{U_{V_s}(x^{(V_s)})\varphi_s(x^{(s)})^{-1} \sum_{y_{W_s}} U_{\{s\}\cup W_s}(x^{(s)} \wedge y^{(W_s)})}{\varphi_s(x^{(s)})^{-1}\left(\sum_{y^{(V_s^*)}} U_{V_s}(x^{(s)} \wedge y^{(V_s^*)})\right)\left(\sum_{y^{(W_s)}} U_{\{s\}\cup W_s}(x^{(s)} \wedge y^{(W_s)})\right)}$$

$$= \frac{U_{V_s}(x^{(V_s)})}{Z_s(x^{(s)})}$$

Now, if $t_1, \ldots, t_n$ are the children of $s$, we have

$$U_{V_s}(x^{(V_s)}) = \varphi_s(x^{(s)}) \prod_{k=1}^{n} \varphi_{st_k}(x^{(s)}, x^{(t_k)}) \prod_{k=1}^{n} U_{V_{t_k}}(x^{(V_{t_k})}),$$

so that

$$\mathbb{P}(X^{(t_k)} = x^{(t_k)}, k = 1, \ldots, n \mid X^{(s)} = x^{(s)})$$

$$= \frac{1}{Z_s(x^{(s)})} \sum_{y^{(V_{t_k}^*)}, k=1,\ldots,n} \varphi_s(x^{(s)}) \prod_{k=1}^{n} \varphi_{st_k}(x^{(s)}, x^{(t_k)}) \prod_{k=1}^{n} U_{V_{t_k}}(x^{(t_k)} \wedge y^{(V_{t_k}^*)})$$

$$= \frac{\varphi_s(x^{(s)}) \prod_{k=1}^{n} \varphi_{st_k}(x^{(s)}, x^{(t_k)}) \prod_{k=1}^{n} Z_{t_k}(x^{(t_k)})}{Z_s(x^{(s)})}$$

This implies that the transition probability needed for the tree model, $p_{st_1}(x^{(s)}, x^{(t_1)})$, must be proportional to $\varphi_{st_1}(x^{(s)}, x^{(t_1)})Z_{t_1}(x^{(t_1)})$ which proves the lemma. ∎

This lemma reduces the computation of the transition probabilities to the computation of $Z_s(x^{(s)})$, for $s \in V$. This can be done efficiently, going upward in the tree (from terminal vertexes to the root). Indeed, if $s$ is terminal, then $V_s = \{s\}$ and $Z_s(x^{(s)}) = \varphi_s(x^{(s)})$. Now, if $s$ is non-terminal and $t_1, \ldots, t_n$ are its children, then, it is easy to see that

$$Z_s(x^{(s)}) = \varphi_s(x^{(s)}) \sum_{x^{(t_1)} \in F_{t_1}, \ldots, x^{(t_n)} \in F_{t_n}} \prod_{k=1}^{n} \varphi_{st_k}(x^{(s)}, x^{(t_k)})Z_{t_k}(x^{(t_k)})$$

$$= \varphi_s(x^{(s)}) \prod_{k=1}^{n} \left( \sum_{x^{(t_k)} \in F_{t_k}} \varphi_{st_k}(x^{(s)}, x^{(t_k)})Z_{t_k}(x^{(t_k)}) \right) \qquad (14.8)$$

So, $Z_s(x^{(s)})$ can be easily computed once the $Z_t(x^{(t)})$'s are known for the children of $s$.

Equations (14.6) to (14.8) therefore provide the necessary relations in order to compute the singleton and edge marginal probabilities on the tree. It is important to note that these relations are valid for any tree structure consistent with the acyclic graph we started with. We now rephrase them with notation that only depend on this graph and not on the selected orientation.

Let $\{s,t\}$ be an edge in $E$. Then $s$ separates the graph $G \setminus \{s\}$ into two components. Let $V_{st}$ be the component that contains $t$, and $V_{st}^* = V_{st} \setminus t$. Define

$$Z_{st}(x_t) = \sum_{y^{(V_{st}^*)} \in \mathcal{F}(V_{st}^*)} U_{V_{st}}(x^{(t)} \wedge y^{(V_{st}^*)}).$$

This $Z_{st}$ coincides with the previously introduced $Z_t$, computed with any tree in which the edge $\{s,t\}$ is oriented from $s$ to $t$. Equation (14.8) can be rewritten with this new notation in the form:

$$Z_{st}(x^{(t)}) = \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} \left( \sum_{x^{(t')} \in F_{t'}} \varphi_{tt'}(x^{(t)}, x^{(t')}) Z_{tt'}(x^{(t')}) \right). \tag{14.9}$$

This equation is usually written in terms of "messages" defined by

$$m_{ts}(x^{(s)}) = \sum_{x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) Z_{st}(x^{(t)})$$

which yields

$$Z_{st}(x^{(t)}) = \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)})$$

and the message consistency relation

$$m_{ts}(x^{(s)}) = \sum_{x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)}). \tag{14.10}$$

Also, because one can start building a tree from $G^\flat$ using any vertex as a root, (14.6) is valid for any $s \in V$, in the form (applying (14.8) to the root)

$$\pi_s(x^{(s)}) = \frac{1}{\zeta_s} \varphi_s(x^{(s)}) \prod_{t \in \mathcal{V}_s} m_{ts}(x^{(s)}) \tag{14.11}$$

where $\zeta_s$ is chosen to ensure that the sum of probabilities is 1. (In fact, looking at lemma 14.1, we have $Z_s = Z$, independent of $s$.)

Similarly, (14.7) can be written

$$p_{st}(x^{(s)}, x^{(t)}) = m_{ts}(x^{(s)})^{-1} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)}) \tag{14.12}$$

which provides the edge transition probabilities. Combining this with (14.11), we get the edge marginal probabilities:

$$\pi_{st}(x^{(s)}, x^{(t)}) = \frac{1}{\zeta}\varphi_{st}(x^{(s)}, x^{(t)})\varphi_s(x^{(s)})\varphi_t(x^{(t)}) \prod_{t'\in\mathcal{V}_t\backslash\{s\}} m_{t't}(x^{(t)}) \prod_{s'\in\mathcal{V}_s\backslash\{t\}} m_{s's}(x^{(s)}). \quad (14.13)$$

**Remark 14.2** We can modify (14.10) by multiplying the right-hand side by an arbitrary constant $q_{ts}$ without changing the resulting estimation of probabilities: this only multiplies the messages by a constant, which cancels after normalization. This remark can be useful in particular to avoid numerical overflow; one can, for example, define $q_{ts} = 1/\sum_{x_s\in F_s} m_{ts}(x_s)$ so that the messages always sum to 1. This is also useful when applying belief propagation (see next section) to loopy networks, for which (14.10) may diverge while the normalized version converges.           ◆

The following summarizes this message passing algorithm.

---

**Algorithm 14.4 (Belief propagation on acyclic graphs)**
Given a family of interactions $\varphi_s : F_s \to [0, +\infty)$, $\varphi_{st} : F_s \times F_t \to [0, +\infty)$,

(1) Initialize functions (messages) $m_{ts} : F_s \to \mathbb{R}$, e.g., taking $m_{ts}(x^{(s)}) = 1/|F_s|$.

(2) Compute unnormalized messages

$$\tilde{m}_{ts}(\cdot) = \sum_{x^{(t)}\in F_t} \varphi_{st}(\cdot, x^{(t)})\varphi_t(x^{(t)}) \prod_{t'\in\mathcal{V}_t\backslash\{s\}} m_{t't}(x^{(t)})$$

and let $m_{ts}(\cdot) = q_{ts}\tilde{m}_{ts}(\cdot)$, for some choice of constant $q_{ts}$, which must be a fixed function of $\tilde{m}_{ts}(\cdot)$, such as

$$q_{ts} = \left( \sum_{x^{(s)}\in F_s} \tilde{m}_{ts}(x^{(s)}) \right)^{-1}.$$

(3) Stop the algorithm when the messages stabilize (which happens after a finite number of updates). Compute the edge marginal distributions using (14.13).

---

It should be clear, from the previous analysis that messages stabilize in finite time, starting from the outskirts of the acyclic graph. Indeed, messages starting from a terminal $t$ (a vertex with only one neighbor) are automatically set to their correct value in (14.10),

$$m_{ts}(x_s) = \sum_{x_t\in F_t} \varphi_{st}(x_s, x_t)\varphi_t(x_t),$$

at the first update. These values then propagate to provide messages that satisfy (14.10) starting from the next-to-terminal vertexes (those that have only one neighbor left when the terminals are removed) and so on.

## 14.3 Belief propagation and free energy approximation

### 14.3.1 BP stationarity

It is possible to run Algorithm 14.4 on graphs that are not acyclic, since nothing in its formulation requires this property. However, while the method stabilizes in finite time for acyclic graphs, this property, or even the convergence of the messages is not guaranteed for general, loopy, graphs. Convergence, however, has been observed in a large number of applications, sometimes with very good approximations of the true marginal distributions.

We will refer to stable solutions of Algorithm 14.4 as BP-stationary points, as formally stated in the next definition, which allows for a possible normalization of messages, which is particularly important with loopy networks.

**Definition 14.3** *Let $G = (V, E)$ be an undirected graph and $\Phi = (\varphi_{st}, \{s, t\} \in E, \varphi_s, s \in V)$ a consistent family of pair interactions. We say that a family of joint probability distributions $(\pi'_{st}, \{s, t\} \in E)$ is BP-stationary for $(G, \Phi)$ if there exists messages $x_t \in F_t \mapsto m_{st}(x_t)$, constants $\zeta_{st}$ for $t \sim s$ and $\alpha_s$ for $s \in V$ satisfying*

$$m_{ts}(x^{(s)}) = \frac{\alpha_s}{\zeta_{ts}} \sum_{x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)}) \qquad (14.14)$$

*such that*

$$\pi'_{st}(x^{(s)}, x^{(t)}) = \frac{1}{\zeta_{st}} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_s(x^{(s)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)}) \prod_{s' \in \mathcal{V}_s \setminus \{t\}} m_{s's}(x^{(s)}). \quad (14.15)$$

There is no loss of generality in the specific form chosen for the normalizing constants in (14.14) and (14.15), in the sense that, if the messages satisfy (14.15) and

$$m_{ts}(x^{(s)}) = q_{ts} \sum_{x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)})$$

for some constants $q_{ts}$, then

$$
\begin{aligned}
\zeta_{st} &= \sum_{x^{(s)} \in F_s, x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_s(x^{(s)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)}) \prod_{s' \in \mathcal{V}_s \setminus \{t\}} m_{s's}(x^{(s)}) \\
&= \frac{1}{q_{ts}} \sum_{x^{(s)} \in F_s} \varphi_s(x^{(s)}) \prod_{s' \in \mathcal{V}_s} m_{s's}(x^{(s)})
\end{aligned}
$$

so that $\zeta_{st} q_{ts}$ (which has been denoted $\alpha_s$) does not depend on $t$. Of course, the relevant questions regarding BP-stationarity is whether the collection of pairwise probability $\pi'_{st}$ exists, how to compute them, and whether $\pi'_{st}(x^{(s)}, x^{(t)})$ provides a good

approximation of the marginals of the probability distribution $\pi$ that is associated to $\Phi$, namely

$$\pi(x) = \frac{1}{Z} \prod_{s \in V} \varphi_s(x^{(s)}) \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}).$$

A reassuring statement for BP-stationarity is that it is not affected when the functions in $\Phi$ are multiplied by constants, which does not affect the underlying probability $\pi$. This is stated in the next proposition.

**Proposition 14.4** *Let $\Phi$ be as above a family of edge and vertex interactions. Let $c_{st}, \{s, t\} \in E, c_s, s \in V$ be families of positive constants, and define $\tilde{\Phi} = (\tilde{\varphi}_{st}, \tilde{\varphi}_s)$ by $\tilde{\varphi}_{st} = c_{st}\varphi_{st}$ and $\tilde{\varphi}_s = c_s\varphi_s$. Then,*

$$\pi' \text{ is BP-stationary for } (G, \Phi) \Leftrightarrow \pi' \text{ is BP-stationary for } (G, \tilde{\Phi}).$$

PROOF Indeed, if (14.14) and (14.15) are true for $(G, \Phi)$, it suffices to replace $\alpha_s$ by $\alpha_s c_s$ and $\zeta_{st}$ by $\zeta_{st} c_{st} c_t$ to obtain (14.14) and (14.15) for $(G, \tilde{\Phi})$. ∎

It is also important to notice that, if $G$ is acyclic, definition 14.3 is no more general than the message-passing rule we had considered earlier. More precisely, we have (see remark 14.2),

**Proposition 14.5** *Let $G = (V, E)$ be undirected acyclic and $\Phi = (\varphi_{st}, \{s, t\} \in E, \varphi_s, s \in V)$ a consistent family of pair interactions. Then, the only BP-stationary distributions are the marginals of the distribution $\pi$ associated to $\Phi$.*

### 14.3.2 Free-energy approximations

A partial justification of the good behavior of BP with general graphs has been provided in terms of a quantity introduced in statistical mechanics, called the Bethe free energy. We let $G = (V, E)$ be an undirected graph and assume that a consistent family of pair interactions is given (denoted $\Phi = (\varphi_s, s \in V, \varphi_{st}, \{s, t\} \in E)$) and consider the associated distribution, $\pi$, on $\mathcal{F}(V)$, given by

$$\pi(x) = \frac{1}{Z} \prod_{s \in V} \varphi_s(x^{(s)}) \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}). \tag{14.16}$$

It will also be convenient to use the function

$$\psi_{st}(x^{(s)}, x^{(t)}) = \varphi_s(x^{(s)})\varphi_t(x^{(t)})\varphi_{st}(x^{(s)}, x^{(t)})$$

such that

$$\pi(x) = \frac{1}{Z} \prod_{s \in V} \varphi_s(x^{(s)})^{1-|\mathcal{V}_s|} \prod_{\{s,t\} \in E} \psi_{st}(x^{(s)}, x^{(t)}). \tag{14.17}$$

We will consider approximations $\pi'$ of $\pi$ that minimize the Kullback-Leibler divergence, $KL(\pi'\|\pi)$ (see (4.3)), subject to some constraints. We can write

$$
\begin{aligned}
KL(\pi'\|\pi) &= -E_{\pi'}(\log \pi) - H(\pi')\\
&= -\log Z - \sum_{s \in V}(1 - |\mathcal{V}_s|)E_{\pi'}(\log \varphi_s) - \sum_{\{s,t\} \in E} E_{\pi'}(\log \psi_{st}) - \mathcal{H}(\pi')
\end{aligned}
$$

(where $\mathcal{H}(\pi')$ is the entropy of $\pi'$). Introduce the one- and two-dimensional marginals of $\pi'$, denoted $\pi'_s$ ad $\pi'_{st}$. Then

$$
\begin{aligned}
KL(\pi'\|\pi) = -\log Z - \sum_{s \in V}(1 - |\mathcal{V}_s|)E_{\pi'}(\log \frac{\varphi_s}{\pi'_s}) - \sum_{\{s,t\} \in E} E_{\pi'}(\log \frac{\psi_{st}}{\pi'_{st}})\\
+ \sum_{s \in V}(1 - |\mathcal{V}_s|)\mathcal{H}(\pi'_s) + \sum_{\{s,t\} \in E} \mathcal{H}(\pi'_{st}) - \mathcal{H}(\pi').
\end{aligned}
$$

The Bethe free energy is the function $\mathbb{F}_\beta$ defined by

$$
\mathbb{F}_\beta(\pi') = -\sum_{s \in V}(1 - |\mathcal{V}_s|)E_{\pi'}(\log \frac{\varphi_s}{\pi'_s}) - \sum_{\{s,t\} \in E} E_{\pi'}(\log \frac{\psi_{st}}{\pi'_{st}}); \tag{14.18}
$$

so that

$$
KL(\pi'\|\pi) = \mathbb{F}_\beta(\pi') - \log Z + \Delta_G(\pi')
$$

with

$$
\Delta_G(\pi') = \sum_{s \in V}(1 - |\mathcal{V}_s|)\mathcal{H}(\pi'_s) + \sum_{\{s,t\} \in E} \mathcal{H}(\pi'_{st}) - \mathcal{H}(\pi').
$$

Using this computation, one can consider the approximation problem: find $\hat{\pi}'$ that minimizes $KL(\pi'\|\pi)$ over a class of distributions $\pi'$ for which the computation of the first and second order marginals is easy. This problem has an explicit solution when the distribution $\pi'$ is such that all variables are independent, leading to what is called the *mean-field approximation* of $\pi$. Indeed, in this case, we have

$$
\Delta_G(\pi') = \sum_{\{s,t\} \in G}(\mathcal{H}(\pi'_s) + \mathcal{H}(\pi'_t)) + \sum_{s \in S}(1 - |\mathcal{V}_s|)\mathcal{H}(\pi'_s) - \sum_{s \in S}\mathcal{H}(\pi'_s) = 0
$$

and

$$
\mathbb{F}_\beta(\pi') = -\sum_{s \in V}(1 - |\mathcal{V}_s|)E_{\pi'}(\log \frac{\varphi_s}{\pi'_s}) - \sum_{\{s,t\} \in E} E_{\pi'}(\log \frac{\psi_{st}}{\pi'_s\pi'_t}).
$$

$\mathbb{F}_\beta$ must be minimized with respect to the variables $\pi'_s(x^{(s)}), s \in S, x_s \in F_S$ subject to the constraints $\sum_{x_s \in F_s} \pi'_s(x^{(s)}) = 1$. The corresponding necessary optimality conditions equations provide the mean-field consistency equations, described in the following definition.

**Proposition 14.6** *A local minimum of $\mathbb{F}_\beta(\pi')$ over all probability distributions $\pi'$ of the form*

$$\pi'(x) = \prod_{s \in V} \pi'_s(x^{(s)})$$

*must satisfy the mean field consistency equations:*

$$\pi_s(x^{(s)}) = \frac{1}{Z_s} \varphi_s(x^{(s)})^{1-|\mathcal{V}_s|} \prod_{t \sim s} \exp\left(E_{\pi_t}(\log \psi_{st}(x^{(.)}.))\right). \tag{14.19}$$

Proof Since all constraints are affine, we can use Lagrange multipliers, denoted $(\lambda_s, s \in S)$ for each of the constraints, to obtain necessary conditions for a minimizer, yielding

$$\frac{\partial \mathbb{F}_\beta}{\partial \pi_s(x_s)} - \lambda_s = 0, \quad s \in S, x_s \in F_s.$$

This gives:

$$-(1 - |\mathcal{V}_s|)\left(\log \frac{\varphi_s(x_s)}{\pi_s(x_s)} - 1\right) - \sum_{t \sim s} \sum_{x_t \in F_t} \left(\log \frac{\psi_{st}(x_s, x_t)}{\pi_s(x_s)\pi_t(x_t)} - 1\right)\pi_t(x_t) = \lambda_s.$$

Solving this with respect to $\pi_s(x_s)$ and regrouping all constant terms (independent from $x_s$) in the normalizing constant $Z_s$ yields (14.19). ∎

The mean field consistency equations can be solved using a root-finding algorithm or by directly solving the minimization problem. We will retrieve this method, with more details, in our discussion of variational approximations in chapter 16.

In the particular case in which $G$ is acyclic and the approximation is made by $G$-Markov processes, the Kullback-Leibler distance is minimized with $\pi' = \pi$ (since $\pi$ belongs to the approximating class). A slightly non-trivial remark is that $\pi$ is optimal also for the minimization of the Bethe free energy $F_\beta$, because this energy coincides, up to the constant term $\log Z$, with the Kullback-Leibler divergence, as proved by the following proposition.

**Proposition 14.7** *If $G$ is acyclic and $\pi'$ is $G$-Markov, then $\Delta_G(\pi') = 0$.*

This proposition is a consequence of the following lemma that has its own interest:

**Lemma 14.8** *If $G$ is acyclic and $\pi$ is a $G$-Markov distribution, then*

$$\pi(x) = \prod_{s \in V} \pi_s(x^{(s)})^{1-|\mathcal{V}_s|} \prod_{\{s,t\} \in E} \pi_{st}(x^{(s)}, x^{(t)}). \tag{14.20}$$

PROOF (OF LEMMA 14.8) We know that, if $\tilde{G} = (V, \tilde{E})$ is a tree such that $\tilde{G}^\flat = G$, we have, letting $s_0$ be the root in $\tilde{G}$

$$
\begin{aligned}
\pi(x) &= \pi_{s_0}(x^{(s_0)}) \prod_{(s,t)\in\tilde{E}} p_{st}(x^{(s)}, x^{(t)}) \\
&= \pi_{s_0}(x^{(s_0)}) \prod_{(s,t)\in\tilde{E}} (\pi_{st}(x^{(s)}, x^{(t)}) \pi(x^{(s)})^{-1}).
\end{aligned}
$$

Each vertex $s$ in $V$ has $|\mathcal{V}_s| - 1$ children in $\tilde{G}$, except $s_0$ which has $|\mathcal{V}_{s_0}|$ children. Using this, we get

$$
\begin{aligned}
\pi(x) &= \pi_{s_0}(x^{(s_0)}) \pi_{s_0}(x^{(s_0)})^{-|\mathcal{V}_{s_0}|} \prod_{s\in V\setminus\{s_0\}} \pi_s(x^{(s)})^{1-|\mathcal{V}_s|} \prod_{(s,t)\in\tilde{E}} \pi_{st}(x^{(s)}, x^{(t)}) \\
&= \prod_{s\in V} \pi_s(x^{(s)})^{1-|\mathcal{V}_s|} \prod_{\{s,t\}\in E} \pi_{st}(x^{(s)}, x^{(t)}).
\end{aligned}
$$

PROOF (OF PROPOSITION 14.7) If $\pi'$ is given by (14.20), then

$$
\begin{aligned}
H(\pi') &= -E_{\pi'} \log \pi' \\
&= -\sum_{s\in V}(1-|\mathcal{V}_s|) E_{\pi'} \log \pi'_s - \sum_{\{s,t\}\in E} E_{\pi'} \log \pi'_{st} \\
&= \sum_{s\in V}(1-|\mathcal{V}_s|) H(\pi'_s) + \sum_{\{s,t\}\in E} H(\pi'_{st})
\end{aligned}
$$

which proves that $\Delta_G(\pi') = 0$. ∎

In view of this, it is tempting to "generalize" the mean field optimization procedure and minimize $\mathbb{F}_\beta(\pi')$ over all possible consistent singletons and pair marginals ($\pi'_s$ and $\pi'_{st}$), then use the optimal ones as an approximation of $\pi_s$ and $\pi_{st}$. What we have just proved is that this procedure provides the exact expression of the marginals when $G$ is acyclic. For loopy graphs, however, it is not justified, and is at best an approximation. A very interesting fact is that this procedure provides the same consistency equations as belief propagation. To see this, we first start with the characterization of minimizers of $\mathbb{F}_\beta$.

**Proposition 14.9** *Let $G = (V, E)$ be an undirected graph and $\pi$ be given by (14.16). Consider the problem of minimizing the Bethe free energy $\mathbb{F}_\beta$ in (14.18) with respect to all possible choices of probability distributions $(\pi'_{st}, \{s, t\} \in E)$, $(\pi'_s, s \in V)$ with the constraints*

$$
\pi'_s(x^{(s)}) = \sum_{x^{(t)}\in F_t} \pi'_{st}(x^{(s)}, x^{(t)}), \forall x^{(s)} \in F_s \text{ and } t \sim s.
$$

*Then a local minimum of this problem must take the form*

$$\pi'_{st}(x^{(s)}, x^{(t)}) = \frac{1}{Z_{st}} \psi_{st}(x^{(s)}, x^{(t)}) \mu_{st}(x^{(t)}) \mu_{ts}(x^{(s)}) \tag{14.21}$$

*where the functions $\mu_{st} : F_t \to [0, +\infty)$ are defined for all $(s,t)$ such that $\{s,t\} \in E$ and satisfy the consistency conditions:*

$$\mu_{ts}(x^{(s)})^{-(|\mathcal{V}_s|-1)} \prod_{s' \sim s} \mu_{s's}(x^{(s)}) = \left( \frac{e}{Z_{st}} \sum_{x^{(t)} \in F_t} \psi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \mu_{st}(x^{(t)}) \right)^{|\mathcal{V}_s|-1}. \tag{14.22}$$

PROOF We introduce Lagrange multipliers: $\lambda_{ts}(x^{(s)})$ for the constraint

$$\pi'_s(x^{(s)}) = \sum_{x^{(t)} \in F_t} \pi'_{st}(x^{(s)}, x^{(t)})$$

and $\gamma_{st}$ for

$$\sum_{x^{(s)}, x^{(t)}} \pi'_{st}(x^{(s)}, x^{(t)}) = 1,$$

which covers all constraints associated to the minimization problem. The associated Lagrangian is

$$\mathbb{F}_\beta(\pi') - \sum_{s \in V} \sum_{x^{(s)} \in F_s} \sum_{t \sim s} \lambda_{ts}(x^{(s)}) \left( \sum_{x^{(t)} \in F_t} \pi'_{st}(x^{(s)}, x^{(t)}) - \pi'_s(x^{(s)}) \right)$$

$$- \sum_{\{s,t\} \in E} \gamma_{st} \left( \sum_{x^{(s)} \in F_s, x^{(t)} \in F_t} \pi'_{st}(x^{(s)}, x^{(t)}) - 1 \right).$$

The derivative with respect to $\pi'_{st}(x^{(s)}, x^{(t)})$ yields the condition

$$\log \pi'_{st}(x^{(s)}, x^{(t)}) - \log \psi_{st}(x^{(s)}, x^{(t)}) + 1 - \lambda_{ts}(x^{(s)}) - \lambda_{st}(x^{(t)}) - \gamma_{st} = 0.$$

which implies

$$\pi'_{st}(x^{(s)}, x^{(t)}) = \varphi_{st}(x^{(s)}, x^{(t)}) \exp(\gamma_{st} - 1) \exp(\lambda_{ts}(x^{(s)}) + \lambda_{st}(x^{(t)})).$$

We let $Z_{st} = \exp(1 - \gamma_{st})$, with $\gamma_{st}$ chosen so that $\pi'_{st}$ is a probability. The derivative with respect to $\pi'_s(x^{(s)})$ gives

$$(1 - |\mathcal{V}_s|)(\log \pi'_s(x^{(s)}) - \log \varphi_s(x^{(s)}) + 1) + \sum_{t \sim s} \lambda_{ts}(x^{(s)}) = 0.$$

Combining this with the expression just obtained for $\pi'_{st}$, we get, for $t \sim s$,

$$(1 - |\mathcal{V}_s|)\log \sum_{x^{(t)} \in F_t} \psi_{st}(x^{(s)}, x^{(t)})e^{\lambda_{st}(x^{(t)})} + (1 - |\mathcal{V}_s|)\lambda_{ts}(x^{(s)})$$

$$+ (1 - |\mathcal{V}_s|)(1 - \log Z_{st} - \log \varphi_s(x^{(s)})) + \sum_{s' \sim s} \lambda_{s's}(x^{(s)}) = 0,$$

which gives (14.22) with $\mu_{st} = \exp(\lambda_{st})$. ∎

A family $\pi'_{st}$ satisfying conditions (14.21) and (14.22) of proposition 14.9 will be called Bethe-consistent. A very interesting remark states that Bethe-consistency is equivalent to BP-stationarity, as stated below.

**Proposition 14.10** *Let $G = (V, E)$ be an undirected graph and $\Phi = (\varphi_{st}, \{s, t\} \in E, \varphi_s, s \in V)$ a consistent family of pair interactions. Then a family $\pi'$ of joint probability distributions is BP-stationary if and only if it is Bethe-consistent.*

PROOF First assume that $\pi'$ is BP-stationary with messages $m_{st}$, so that (14.14) and (14.15) are satisfied. Take

$$\mu_{st} = a_t \prod_{t' \in \mathcal{V}_t, t' \neq s} m_{t't}(x^{(t)})$$

for some constant $a_t$ that will be determined later. Then, the left-hand side of (14.22) is

$$\mu_{ts}(x^{(s)})^{-(|\mathcal{V}_s|-1)} \prod_{s' \in \mathcal{V}_s} \mu_{s's}(x^{(s)}) = a_s \left( \prod_{s' \in \mathcal{V}_s, s' \neq t} m_{s's}(x^{(s)}) \right)^{-(|\mathcal{V}_s|-1)} \prod_{s' \in \mathcal{V}_s} \prod_{s'' \in \mathcal{V}_s, s'' \neq s'} m_{s''s}(x^{(s)})$$

$$= a_s m_{ts}(x^{(s)})^{|\mathcal{V}_s|-1}.$$

The right-hand side is equal to (using (14.14))

$$\left( \frac{ea_t \zeta_{st}}{Z_{st} \alpha_s} m_{ts}(x^{(s)}) \right)^{|\mathcal{V}_s|-1},$$

so that we need to have

$$a_s = \left( \frac{ea_t \zeta_{st}}{Z_{st} \alpha_s} \right)^{|\mathcal{V}_s|-1}.$$

We also need

$$Z_{st} = \sum_{x^{(s)}, x^{(t)}} \psi_{st}(x^{(s)}, x^{(t)})\mu_{st}(x^{(t)})\mu_{ts}(x^{(s)}) = a_s a_t \zeta_{st}.$$

Solving these equations, we find that (14.21) and (14.22) are satisfied with

$$\begin{cases} a_s = (e/\alpha_s)^{(|\mathcal{V}_s|-1)/|\mathcal{V}_s|} \\ Z_{st} = \zeta_{st} a_s a_t \end{cases}$$

which proves that $\pi'$ is Bethe-consistent.

Conversely, take a Bethe-consistent $\pi'$, and $\mu_{st}, Z_{st}$ satisfying (14.21) and (14.22). For $s$ such that $|\mathcal{V}_s| > 1$, define, for $t \in \mathcal{V}_s$,

$$m_{ts}(x^{(s)}) = \mu_{ts}(x^{(s)})^{-1} \prod_{s' \sim s} \mu_{s's}(x^{(s)})^{1/(|\mathcal{V}_s|-1)}. \qquad (14.23)$$

Define also, for $|\mathcal{V}_s| > 1$,

$$\rho_{ts}(x^{(s)}) = \prod_{s' \in \mathcal{V}_s, s' \neq t} m_{s's}(x^{(s)}).$$

(If $|\mathcal{V}_s| = 1$, take $\rho_{ts} \equiv 1$.) Using (14.23), we find $\rho_{ts} = \mu_{ts}$ when $|\mathcal{V}_s| > 1$, and this identity is still valid when $|\mathcal{V}_s| = 1$, since in this case, (14.22) implies that $\mu_{ts}(x^{(s)}) = 1$.

We need to find constants $\alpha_t$ and $\zeta_{st}$ such that (14.14) and (14.15) are satisfied. But (14.15) implies

$$\zeta_{ts} = \sum_{x_t, x_s} \psi_{st}(x^{(s)}, x^{(t)}) \rho_{st}(x^{(t)}) \rho_{ts}(x^{(s)})$$

and (14.21) implies $\zeta_{ts} = Z_{ts}$.

We now consider (14.14), which requires

$$m_{ts}(x^{(s)}) = \frac{\alpha_s}{\zeta_{st}} \sum_{x^{(t)}} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \rho_{st}(x^{(t)}).$$

It is now easy to see that this identity to the power $|\mathcal{V}_s| - 1$ coincides with (14.22) as soon as one takes $\alpha_s = e$.                                                                ∎

## 14.4   Computing the most likely configuration

We now address the problem of finding a configuration that maximizes $\pi(x)$ (mode determination). This problem turns out to be very similar to the computation of marginals, that we have considered so far, and we will obtain similar algorithms.

Assume that $G$ is undirected and acyclic and that $\pi$ can be written as

$$\pi(x) = \frac{1}{Z} \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}) \prod_{s \in V} \varphi_s(x^{(s)}).$$

Maximizing $\pi(x)$ is equivalent to maximizing

$$U(x) = \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}) \prod_{s \in V} \varphi_s(x^{(s)}). \qquad (14.24)$$

Assume that a root has been chosen in $G$, with the resulting edge orientation yielding a tree $\tilde{G} = (V, \tilde{E})$ such that $\tilde{G}^\flat = G$. We partially order the vertexes according to $\tilde{G}$, writing $s \leq t$ if there exists a path from $s$ to $t$ in $\tilde{G}$ ($s$ is an ancestor of $t$). Let $V_s^+$ contain all $t \in V$ with $t \geq s$, and define

$$U_s(x^{(V_s^+)}) = \prod_{\{t,u\} \in E_{V_s^+}} \varphi_{tu}(x^{(t)}, x^{(u)}) \prod_{t>s} \varphi_t(x^{(t)})$$

and

$$U_s^*(x^{(s)}) = \max\left\{U_s(y^{(V_s^+)}), y^{(s)} = x^{(s)}\right\}. \tag{14.25}$$

Since we can write

$$U_s(x^{(V_s^+)}) = \prod_{t \in s^+} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) U_t(x^{(V_t^+)}), \tag{14.26}$$

we have

$$
\begin{aligned}
U_s^*(x^{(s)}) &= \max_{x^{(t)}, t \in s^+} \left( \prod_{t \in s^+} \varphi_t(x^{(t)}) \varphi_{st}(x^{(s)}, x^{(t)}) U_t^*(x^{(t)}) \right) \\
&= \prod_{t \in s^+} \max_{x_t \in F_t} (\varphi_t(x^{(t)}) \varphi_{st}(x^{(s)}, x^{(t)}) U_t^*(x^{(t)})).
\end{aligned}
\tag{14.27}
$$

This provides a method to compute $U_s^*(x^{(s)})$ for all $s$, starting with the leaves and progressively updating the parents. (When $s$ is a leaf, $U_s^*(x^{(s)}) = 1$, by definition.)

Once all $U_s^*(x^{(s)})$ have been computed, it is possible to obtain a configuration $x_*$ that maximizes $\pi$. This is because an optimal configuration must satisfy $U_s^*(x_*^{(s)}) = U_s(x_*^{(V_s^+)})$ for all $s \in V$, i.e., $x_*^{(V_s^+ \setminus \{s\})}$ must solve the maximization problem in (14.25). But because of (14.26), we can separate this problem over the children of $s$ and obtain the fact that, it $t \in s^+$,

$$x_*^{(t)} = \operatorname*{argmax}_{x^{(t)}} \left( \varphi_t(x^{(t)}) \varphi_{st}(x_*^{(s)}, x^{(t)}) U_t^*(x^{(t)}) \right).$$

This procedure can be rewritten in a slightly different form using messages similar to the belief propagation algorithm. It $s \in t^+$, define

$$\mu_{st}(x^{(t)}) = \max_{x_s \in F_s} (\varphi_t(x^{(t)}) \varphi_{ts}(x^{(t)}, x^{(s)}) U_s^*(x^{(s)}))$$

and

$$\xi_{st}(x^{(t)}) = \operatorname*{argmax}_{x^{(s)} \in F_s} (\varphi_t(x^{(t)}) \varphi_{ts}(x^{(t)}, x^{(s)}) U_s^*(x^{(s)})).$$

Using section 14.4, we get

$$
\mu_{st}(x^{(t)}) \;=\; \max_{x^{(s)} \in F_s}\left(\varphi_{ts}(x^{(t)}, x^{(s)})\varphi_s(x^{(s)}) \prod_{u \in s^+} \mu_{us}(x^{(s)})\right),
$$

$$
\xi_{st}(x_t) \;=\; \operatorname*{argmax}_{x^{(s)} \in F_s}\left(\varphi_{ts}(x^{(t)}, x^{(s)})\varphi_s(x^{(s)}) \prod_{u \in s^+} \mu_{us}(x^{(s)})\right).
$$

An optimal configuration can now be computed using $x_*^{(t)} = \xi_{ts}(x_*^{(s)})$, with $s \in pa(t)$.

This resulting algorithm therefore first operates upwards on the tree (from leaves to root) to compute the $\mu_{st}$'s and $\xi_{st}$'s, then downwards to compute $x_*$. This is summarized in the following algorithm.

---

**Algorithm 14.5**
A most likely configuration for

$$
\pi(x) = \frac{1}{Z} \prod_{\{s,t\} \in E} \varphi_{st}(x_s, x_t) \prod_{s \in V} \varphi_s(x_s).
$$

can be computed after iterating the following updates, based on any acyclic orientation of $G$:

(1) Compute, from leaves to root:

$$
\mu_{st}(x^{(t)}) = \max_{x^{(s)} \in F_s}\left(\varphi_{ts}(x^{(t)}, x^{(s)})\varphi_s(x^{(s)}) \prod_{u \in s^+} \mu_{us}(x^{(s)})\right)
$$

and $\xi_{st}(x^{(t)}) = \operatorname*{argmax}\limits_{x^{(s)} \in F_s}\left(\varphi_{ts}(x^{(t)}, x^{(s)})\varphi_s(x^{(s)}) \prod\limits_{u \in s^+} \mu_{us}(x^{(s)})\right).$

(2) Compute, from root to leaves: $x_*^{(t)} = \xi_{ts}(x_*^{(s)})$, with $s = pa(t)$.

---

Similar to the computation of marginals, this algorithm can be rewritten in an orientation-independent form. The main remark is that the value of $\mu_{st}(x^{(t)})$ does not depend on the tree orientation, as long as it is chosen such that $s \in t^+$, i.e., the edge $\{s, t\}$ is oriented from $t$ to $s$. This is because such a choice uniquely prescribes the orientation of the edges of the descendants of $s$ for any such tree, and $\mu_{st}$ only depends on this structure. Since the same remark holds for $\xi_{st}$, this provides a definition of these two quantities for any pair $s, t$ such that $\{s, t\} \in E$. The updating rule

now becomes

$$\mu_{st}(x^{(t)}) = \max_{x^{(s)} \in F_s} \left( \varphi_{ts}(x^{(t)}, x^{(s)}) \varphi_s(x^{(s)}) \prod_{u \in \mathcal{V}_s \setminus \{t\}} \mu_{us}(x^{(s)}) \right), \tag{14.28}$$

$$\xi_{st}(x^{(t)}) = \operatorname*{argmax}_{x^{(s)} \in F_s} \left( \varphi_{ts}(x^{(t)}, x^{(s)}) \varphi_s(x^{(s)}) \prod_{u \in \mathcal{V}_s \setminus \{t\}} \mu_{us}(x^{(s)}) \right) \tag{14.29}$$

with $x_*^{(t)} = \xi_{ts}(x_*^{(s)})$ for any pair $s \sim t$. Like with the $m_{ts}$ in the previous section, looping over updating all $\mu_{ts}$ in any order will finally stabilize to their correct values, although, if an orientation is given, going from leaves to roots is obviously more efficient.

The previous analysis is not valid for loopy graphs but section 14.4 and section 14.4 provide well defined iterations when $G$ is an arbitrary undirected graph, and can therefore be used as such, without any guaranteed behavior.

## 14.5  General sum-prod and max-prod algorithms

### 14.5.1  Factor graphs

The expressions we obtained for message updating with belief propagation and with mode determination respectively took the form

$$m_{ts}(x^{(s)}) \leftarrow \sum_{x^{(t)} \in F_t} \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} m_{t't}(x^{(t)})$$

and

$$\mu_{ts}(x^{(s)}) \leftarrow \max_{x^{(t)} \in F_t} \left( \varphi_{st}(x^{(s)}, x^{(t)}) \varphi_t(x^{(t)}) \prod_{t' \in \mathcal{V}_t \setminus \{s\}} \mu_{t't}(x^{(t)}) \right).$$

They first one is often referred to as the "sum-prod" update rule, and the second as the "max-prod". In our construction, the sum-prod algorithm provided us with a method computing

$$\sigma_s(x^{(s)}) = \sum_{y^{(V \setminus \{s\})}} U(x^{(s)} \wedge y^{(V \setminus \{s\})})$$

with

$$U(x) = \prod_s \varphi_s(x^{(s)}) \prod_{\{s,t\} \in E} \varphi_{st}(x^{(s)}, x^{(t)}).$$

Indeed, we have, according to (14.11)

$$\sigma_s(x^{(s)}) = \varphi_s(x^{(s)}) \prod_{t \in \mathcal{V}_s} m_{ts}(x^{(s)}).$$

Similarly, the max-prod algorithm computes

$$\rho_s(x^{(s)}) = \max_{y_{V \setminus \{s\}}} U(x^{(s)} \wedge y^{(V \setminus \{s\})})$$

via the relation

$$\rho_s(x^{(s)}) = \varphi_s(x^{(s)}) \prod_{t \in \mathcal{V}_s} \mu_{ts}(x^{(s)}).$$

We now discuss generalizations of these algorithms to situations in which the function $U$ does not decompose as a product of bivariate functions. More precisely, let $\mathcal{S}$ be a subset of $\mathcal{P}(V)$, and assume the decomposition

$$U(x) = \prod_{C \subset \mathcal{S}} \varphi_C(x_C).$$

The previous algorithms can be generalized using the concept of *factor graphs* associated with the decomposition. The vertexes of this graph are either indexes $s \in V$ or sets $C \in \mathcal{S}$, and the only edges link indexes and sets that contain them. The formal definition is as follows.

**Definition 14.11** *Let $V$ be a finite set of indexes and $\mathcal{S}$ a subset of $\mathcal{P}(V)$. The factor graph associated to $V$ and $\mathcal{S}$ is the graph $G = (V \cup \mathcal{S}, E)$, $E$ being constituted of all pairs $\{s, C\}$ with $C \in \mathcal{S}$ and $s \in C$.*

We assign the variable $x^{(s)}$ to a vertex $s \in V$ of the factor graph, and the function $\varphi_C$ to $C \in \mathcal{S}$. With this in mind, the sum-prod and max-prod algorithms are extended to factor graphs as follows.

**Definition 14.12** *Let $G = (V \cup \mathcal{S}, E)$ be a factor graph, with associated functions $\varphi_C(x_C)$. The sum-prod algorithm on $G$ updates messages $m_{sC}(x_s)$ and $m_{Cs}(x_s)$ according to the rules*

$$\begin{cases} m_{sC}(x^{(s)}) \leftarrow \displaystyle\prod_{\tilde{C}, s \in \tilde{C}, \tilde{C} \neq C} m_{\tilde{C}s}(x^{(s)}) \\ m_{Cs}(x^{(s)}) \leftarrow \displaystyle\sum_{y_C : y^{(s)} = x^{(s)}} \varphi_C(y^{(C)}) \prod_{t \in C \setminus \{s\}} m_{tC}(y^{(t)}) \end{cases} \quad (14.30)$$

*Similarly, the max-prod algorithm iterates*

$$
\begin{cases}
\mu_{sC}(x^{(s)}) \leftarrow \displaystyle\prod_{\tilde{C}, s\in\tilde{C}, \tilde{C}\neq C} \mu_{\tilde{C}s}(x^{(s)}) \\[2em]
\mu_{Cs}(x^{(s)}) \leftarrow \displaystyle\max_{y^{(C)}:y^{(s)}=x^{(s)}} \varphi_C(y^{(C)}) \prod_{t\in C\setminus\{s\}} \mu_{tC}(y^{(t)})
\end{cases}
\tag{14.31}
$$

These algorithms reduce to the original ones when only single vertex and pair interactions exist. Let us check this with sum-prod. In this case, the set $\mathcal{S}$ contains all singletons $C = \{s\}$, with associated function $\varphi_s$, and all edges $\{s,t\}$ with associated function $\varphi_{st}$. We have links between $s$ and $\{s\}$ and $s$ and $\{s,t\} \in E$. For singletons, we have

$$
m_{s\{s\}}(x^{(s)}) \leftarrow \prod_{t\sim s} m_{s\{s,t\}}(x^{(s)}) \text{ and } m_{\{s\}s}(x^{(s)}) \leftarrow \varphi_s(x^{(s)}).
$$

For pairs,

$$
m_{s\{s,t\}}(x^{(s)}) \leftarrow \varphi_s(x^{(s)}) \prod_{\tilde{t}\in\mathcal{V}_s\setminus\{t\}} m_{\{s,\tilde{t}\}s}(x^{(s)})
$$

and

$$
m_{\{s,t\}s}(x^{(s)}) \leftarrow \sum_{y^{(t)}} \varphi_{st}(x^{(s)}, y^{(t)}) m_{t\{s,t\}}(y^{(t)})
$$

and, combining the last two assignments, it becomes clear that we retrieve the initial algorithm with $m_{\{s,t\}s}$ taking the role of what we previously denoted $m_{ts}$.

The important question, obviously, is whether the algorithms converge. The following result shows that this is true when the factor graph is acyclic.

**Proposition 14.13** *Let $G = (V \cup \mathcal{S}, E)$ be a factor graph with associated functions $\varphi_C$. Assume that $G$ is acyclic. Then the sum-prod and max-prod algorithms converge in finite time.*

*After convergence, we have $\sigma_s(x^{(s)}) = \prod_{C, s\in C} m_{Cs}(x^{(s)})$ and $\rho_s(x^{(s)}) = \prod_{C, s\in C} \mu_{Cs}(x^{(s)})$.*

PROOF Let us assume that $G$ is connected, which is without loss of generality, since the following argument can be applied to each component of $G$ separately. Since $G$ is acyclic, we can arbitrarily select one of its vertexes as a root to form a tree. This being done, we can see that the messages going upward in the tree (from children to parent) progressively stabilize, starting with leaves. Leaves in the factor graph indeed are either singletons, $C = \{s\}$, or vertexes $s \in V$ that belong to only one set $C \in \mathcal{S}$. In the first case, the algorithm imposes (taking, for example, the sum-prod case) $m_{\{s\}s}(x^{(s)}) = \varphi_s(x^{(s)})$, and in the second case $m_{sC}(x^{(s)}) = 1$. So the messages sent upward by the leaves are set at the first step. Since the messages going from a child to its parents only depend on the messages that it received from its other neighbors

in the acyclic graph, which are its children in the tree, it is clear that all upward messages progressively stabilize until the root is reached. Once this is done, messages propagate downward from each parent to its children. This stabilizes as soon as all incoming messages to the parent are stabilized, since outgoing messages only depend on those. At the end of the upward phase, this is true for the root, which can then send its stable message to its children. These children now have all their incoming messages and can now send their messages to their own children and so on down to the leaves.

We now consider the second statement, proceeding by induction, assuming that the result is true for any smaller graph than the one considered. Let $s_0$ be the selected root, and consider all vertexes $s \neq s_0$ such that there exists $C_s \in \mathcal{S}$ such that $s_0$ and $s$ both belong to $C_s$. Given $s$, there cannot be more than one such $C_s$ since this would create a loop in the graph. For each such $s$, consider the part $G_s$ of $G$ containing all descendants of $s$. Let $V_s$ be the set of vertexes among the descendants of $s$ and $\mathcal{C}_s$ the set of $C$'s below $s$. Define

$$U_s(x^{(V_s)}) = \prod_{C \in \mathcal{C}_s} \varphi_C(x^{(C)}).$$

Since the upward phase of the algorithm does not depend on the ancestors of $s$, the messages incoming to $s$ for the sum-prod algorithm restricted to $G_s$ are the same as with the general algorithm, so that, using the induction hypothesis

$$\sum_{y^{(V_s)}, y^{(s)} = x^{(s)}} U_s(y^{(V_s)}) = \prod_{C \in \mathcal{C}_s, s \in C} m_{Cs}(x^{(s)}) = m_{sC_s}(x^{(s)}).$$

Now let $C_1, \ldots, C_n$ list all the sets in $\mathcal{C}$ that contain $s_0$, which must be non-intersecting (excepted at $\{s_0\}$), again not to create loops. Write

$$C_1 \cup \cdots \cup C_n = \{s_0, s_1, \ldots, s_q\}.$$

Then, we have

$$U(x) = \prod_{j=1}^{n} \varphi_{C_j}(x^{(C_j)}) \prod_{i=1}^{q} U_{s_i}(x^{(V_{s_i})})$$

and letting $S = \bigcup_{j=1}^{n} C_j \setminus \{s_0\}$,

$$
\begin{aligned}
\sigma_{s_0}(x^{(s_0)}) &= \sum_{y^{(V)}:y^{(s_0)}=x^{(s_0)}} \prod_{j=1}^{n} \varphi_{C_j}(y^{(C_j)}) \prod_{i=1}^{q} U_{s_i}(y^{(V_{s_i})}) \\
&= \sum_{y^{()}S:y^{(s_0)}=x^{(s_0)}} \prod_{j=1}^{n} \varphi_{C_j}(y^{(C_j)}) \prod_{i=1}^{q} m_{s_i C_{s_i}}(y^{(s_i)}) \\
&= \prod_{j=1}^{n} \sum_{y^{()}C_j:y^{(s_0)}=x^{(s_0)}} \varphi_{C_j}(y^{(C_j)}) \prod_{s \in C_j \setminus \{s_0\}} m_{s C_s}(y^{(s)}) \\
&= \prod_{j=1}^{n} m_{C_j s_0}(x^{(s_0)})
\end{aligned}
$$

which proves the required result (note that, when factorizing the sum, we have used the fact that the sets $C_j \setminus \{s_0\}$ are non intersecting). An almost identical argument holds for the max-prod algorithm. ∎

**Remark 14.14** Note that these algorithms are not always feasible. For example, it is always possible to represent a function $U$ on $\mathcal{F}(V)$ with the trivial factor graph in which $\mathcal{S} = \{V\}$ and $E$ contains all $\{s, V\}, s \in V$ (using $\varphi_V = U$), but computing $m_{Vs}$ is identical to directly computing $\sigma_s$ with a sum over all configurations on $V \setminus \{s\}$ which grows exponentially. In fact, the complexity of the sum-prod and max-prod algorithms is exponential in the size of the largest $C$ in $\mathcal{S}$ which should therefore remain small. ◆

**Remark 14.15** It is not always possible to decompose a function so that the resulting factor graph is acyclic with small degree (maximum number of edges per vertex). Sum-prod and max-prod can still be used with loopy networks, sometimes with excellent results, but without theoretical support. ◆

**Remark 14.16** One can sometimes transform a given factor graph into an acyclic one by grouping vertexes. Assume that the set $\mathcal{S} \subset \mathcal{P}(V)$ is given. We will say that a partition $\Delta = (D_1, \ldots D_k)$ of $V$ is $\mathcal{S}$-admissible if, for any $C \in \mathcal{S}$ and any $j \in \{1, \ldots, k\}$, one has either $D_j \cap C = \emptyset$ or $D_j \subset C$.

If $\Delta$ is $\mathcal{S}$-admissible, one can define a new factor graph $\tilde{G}$ as follows. We first let $\tilde{V} = \{1, \ldots, k\}$. To define $\tilde{\mathcal{S}} \subset \mathcal{P}(\tilde{V})$ assign to each $C \in \mathcal{S}$ the set $J_C$ of indexes $j$ such that $D_j \subset C$. From the admissibility assumption,

$$
C = \bigcup_{j \in J_C} D_j, \tag{14.32}
$$

so that $C \mapsto J_C$ is one-to-one. Let $\tilde{S} = \{J_C, C \in S\}$. Group variables using $\tilde{x}^{(k)} = x^{(D_k)}$, so that $\tilde{F}_k = \mathcal{F}(D_k)$. Define $\tilde{\Phi} = (\tilde{\varphi}_{\tilde{C}}, \tilde{C} \in \tilde{S})$ by $\tilde{\varphi}_{\tilde{C}} = \varphi_C$ where $C$ is given by (14.32).

In other terms, one groups variables $(x^{(s)}, s \in V)$ into clusters, to create a simpler factor graph, which may be acyclic even if the original one was not. For example, if $V = \{a, b, c, d\}$, $S = \{A, B\}$ with $A = \{a, b, c\}$ and $B = \{b, c, d\}$, then $(A, c, B, b)$ is a cycle in the associated factor graph. If, however, one takes $D_1 = \{a\}$, $D_2 = \{b, c\}$ and $D_3 = \{d\}$, then $(D_1, D_2, D_3)$ is $S$-admissible and the associated factor graph is acyclic. In fact, in such a case, the resulting factor graph, considered as a graph with vertexes given by subsets of $V$, is a special case of a junction tree, which is defined in the next section.◆

### 14.5.2  Junction trees

**Definition 14.17** *Let $V$ be a finite set. A junction tree on $V$ is an undirected acyclic graph $\mathbb{G} = (S, \mathbb{E})$ where $S \subset \mathcal{P}(V)$ is a family of subsets of $V$ that satisfy the following property, called the running intersection constraint: if $C, C' \in S$ and $s \in C \cap C'$, then all sets $C''$ in the (unique) path connecting $C$ and $C'$ in $\mathbb{G}$ must also contain $s$.*

**Remark 14.18** Let us check that the clustered factor graph $\tilde{G}$ defined in remark 14.16 is equivalent to a junction tree when acyclic. Using the same notation, let $\hat{S} = \{D_1, \ldots, D_k\} \cup S$, removing if needed sets $C \in S$ that coincide with one of the $D_j$'s. Place an edge between $D_j$ and $C$ if and only if $D_j \subset C$.

Let $(C_1, D_{i_1}, \ldots, D_{i_{n-1}}, C_n)$ be a path in that graph. Assume that $s \in C_1 \cap C_2$. Let $D_{i_n}$ be the unique $D_j$ that contains $s$. It is such that from the the admissibility assumption, $D_{i_n} \subset C_1$ and $D_{i_n} \subset C_n$, which implies that $(C_1, D_{i_1}, \ldots, C_n, D_{i_n}, C_1)$ is a path in $\tilde{G}$. Since $\tilde{G}$ is acyclic, this path must be a union of folded paths. But it is easy to see that any folded path satisfies the running intersection constraint. (Note that there was no loss of generality in assuming that the path started and ended with a "$C$", since any "$D$" must be contained in the $C$ that follows or precedes it.)              ◆

We now consider a probability distribution written in the form

$$\pi(x) = \frac{1}{Z} \prod_{C \in S} \varphi_C(x^{(C)})$$

and we make the assumption that $S$ can be organized as a junction tree.

Belief propagation can be extended to junction trees. Fixing a root $C_0 \in S$, we first choose an orientation on $\mathbb{G}$, which induces as usual a partial order on $S$. For $C \in S$, define $S_C^+$ as the set of all $B \in S$ such that $B > C$. Define also

$$V_C^+ = \bigcup_{B \in S_C^+} B.$$

We want to compute sums

$$\sigma_C(x^{(C)}) = \sum_{y^{(V \setminus C)}} U(x^{(C)} \wedge y^{(V \setminus C)}),$$

where $U(x) = \prod_{C \in \mathcal{S}} \varphi_C(x^{(C)})$. We have

$$\sigma_C(x^{(C)}) = \sum_{y^{(V \setminus C)}} \varphi_C(x^{(C)}) \prod_{B \in \mathcal{S} \setminus \{C\}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}).$$

Define

$$\sigma_C^+(x^{(C)}) = \sum_{y^{(V_C^+ \setminus C)}} \prod_{B > C} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}).$$

Note that we have $\sigma_{C_0} = \varphi_{C_0} \sigma_{C_0}^+$ at the root. We have the recursion formula

$$
\begin{aligned}
\sigma_C^+(x^{(C)}) &= \sum_{y^{(V_C^+ \setminus C)}} \prod_{C \to B} \left( \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \prod_{B' > B} \varphi_{B'}(x^{(B' \cap C)} \wedge y^{(B' \setminus C)}) \right) \\
&= \prod_{C \to B} \sum_{y^{(B \cup V_B^+ \setminus C)}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \prod_{B' > B} \varphi_{B'}(x^{(B' \cap C)} \wedge y^{(B' \setminus C)}) \\
&= \prod_{C \to B} \sum_{y^{(B \setminus C)}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \sigma_B^+(x^{(B \cap C)} \wedge y^{(B \setminus C)}).
\end{aligned}
$$

The inversion between the sum and product in the second equation above was possible because the sets $B \cup V_B^+ \setminus C$, $C \to B$ are disjoint. Indeed, if there existed $B, B'$ such that $C \to B$ and $C \to B'$, and descendants $C'$ of $B'$ and $C''$ of $B''$ with a nonempty intersection, then this intersection would have to be included in every set in the (non-oriented) path connecting $C'$ and $C''$ in $\mathbb{G}$. Since this path contains $C$, the intersection must also be included in $C$, so that the sets $B \cup V_B^+ \setminus C$, with $C \to B$ are disjoint.

Introduce messages

$$m_B^+(x^{(C)}) = \sum_{y^{(B \setminus C)}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \sigma_B^+(x^{(B \cap C)} \wedge y^{(B \setminus C)})$$

where $C$ is the parent of $B$. Then

$$m_B^+(x^{(C)}) = \sum_{y^{(B \setminus C)}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \prod_{B \to B'} m_{B'}^+(x^{(B \cap C)} \wedge y^{(B \setminus C)})$$

with

$$\sigma_C^+(x^{(C)}) = \prod_{C \to B} m_B^+(x^{(C)})$$

which provides $\sigma_C$ at the root. Reinterpreting this discussion in terms of the undirected graph, we are led to introducing messages $m_{BC}(x^{(C)})$ for $B \sim C$ in $\mathbb{G}$, with the message-passing rule

$$m_{BC}(x^{(C)}) = \sum_{y^{(B \setminus C)}} \varphi_B(x^{(B \cap C)} \wedge y^{(B \setminus C)}) \prod_{B' \sim B, B' \neq C} m_{B'B}(x^{(B \cap C)} \wedge y^{(B \setminus C)}). \tag{14.33}$$

Messages progressively stabilize when applied in $\mathbb{G}$, and at convergence, we have

$$\sigma_C(x^{(C)}) = \varphi_C(x^{(C)}) \prod_{B \sim C} m_{BC}(x^{(C)}). \tag{14.34}$$

Note that the complexity of the junction tree algorithm is exponential in the cardinality of the largest $C \in \mathcal{S}$. This algorithm will therefore be unfeasible if $\mathcal{S}$ contains sets that are too large.

## 14.6  Building junction trees

There is more than one family of set interactions with respect to which a given probability $\pi$ can be decomposed (notice that, unlike in the Hammersley-Clifford Theorem, we do not assume that the interactions are normalized), and not all of them can be organized as a junction tree. One can however extend any given family into a new one on which one can build a junction tree.

**Definition 14.19** *Let $V$ be a set of vertexes, and $\mathcal{S}_0 \subset \mathcal{P}(V)$. We say that a set $\mathcal{S} \subset \mathcal{P}(V)$ is an extension of $\mathcal{S}_0$ if, for any $C_0 \in \mathcal{S}_0$, there exists a $C \in \mathcal{S}$ such that $C_0 \subset C$.*

*A tree $\mathbb{G} = (\mathcal{S}, E)$ is a junction-tree extension of $\mathcal{S}_0$ if $\mathcal{S}$ is an extension of $\mathcal{S}_0$ and $\mathbb{G}$ is a junction tree.*

If $\Phi^0 = (\varphi_C^0, C \in \mathcal{S}_0)$ is a consistent family of set interactions, and $\mathcal{S}$ is an extension of $\mathcal{S}_0$, one can build a new family, $\Phi = (\varphi_C, C \in \mathcal{S})$, of set interactions which yields the same probability distribution, i.e., such that, for all $x \in \mathcal{F}(V)$,

$$\prod_{C \in \mathcal{S}} \varphi_C(x^{(C)}) \propto \prod_{C_0 \in \mathcal{S}_0} \varphi_{C_0}^0(x^{(C_0)}).$$

For this, it suffices to build a mapping say $T : \mathcal{S}_0 \to \mathcal{S}$ such that $C_0 \subset T(C_0)$ for all $C_0 \in \mathcal{S}_0$, which is always possible since $\mathcal{S}$ is an extension of $\mathcal{S}_0$ (for example, arbitrarily order the elements of $\mathcal{S}$ and let $T(\mathcal{S}_0)$ be the first element of $\mathcal{S}$, according to this order, that contains $C_0$). One can then define

$$\varphi_C(x^{(C)}) = \prod_{C_0 : T(C_0) = C} \varphi_{C_0}^0(x^{(C_0)}).$$

Given $\Phi^0$, our goal is to design a junction-tree extension which is as feasible as possible. So, we are not interested by the trivial extension $\mathbb{G} = (V, \emptyset)$, since the resulting junction-tree algorithm is unfeasible as soon as $V$ is large. Theorem 14.24 in the next section will be the first step in the design of an algorithm that computes junction trees on a given graph.

### 14.6.1 Triangulated graphs

**Definition 14.20** *Let $G = (V, E)$ be an undirected graph. Let $(s_1, s_2, \ldots, s_n)$ be a path in $G$. One says that this path has a chord at $s_j$, with $j \in \{2, \ldots, n\}$, if $s_{j-1} \sim s_{j+1}$, and we will refer to $(s_{j-1}, s_j, s_{j+1})$ as a chordal triangle. A path in $G$ is achordal if it has no chord.*

*One says that $G$ is triangulated (or chordal) if it has no achordal loop.*

**Definition 14.21** *The graph $G$ is decomposable if it satisfies the following recursive condition: it is either complete, or there exists disjoint subsets $(A, B, C)$ of $V$ such that*

- *$V = A \cup B \cup C$,*
- *$A$ and $B$ are not empty,*
- *$C$ is clique in $G$, $C$ separates $A$ and $B$,*
- *the restricted graphs, $G_{A \cup C}$ and $G_{B \cup C}$ are decomposable.*

These definitions are in fact equivalent, as stated in the following proposition.

**Proposition 14.22** *An undirected graph is triangulated if and only if it is decomposable*

PROOF To prove the "if" part, we proceed by induction on $n = |V|$. Note that every graph for $n \leq 3$ is both decomposable and triangulated (we leave the verification to the reader). Assume that the statement "decomposable $\Rightarrow$ triangulated" holds for graphs with less than $n$ vertexes, and take $G$ with $n$ vertexes. Assume that $G$ is decomposable. If it is complete, it is obviously triangulated. Otherwise, there exists $A, B, C$ such that $V = A \cup B \cup C$, with $A$ and $B$ non-empty such that $G_{A \cup C}$ and $G_{B \cup C}$ are decomposable, hence triangulated from the induction hypothesis, and such that $C$ is a clique which separates $A$ and $B$. Assume that $\gamma$ is an achordal loop in $G$. Since it cannot be included in $A \cup C$ or $B \cup C$, $\gamma$ must go from $A$ to $B$ and back, which implies that it passes at least twice in $C$. Since $C$ is complete, the original loop can be shortcut to form subloops in $A \cup C$ and $B \cup C$. If one of (or both) these loops has cardinality 3, this would provide $\gamma$ with a chord, which contradicts the assumption. Otherwise, the following lemma also provides a contradiction, since one of the two chords that it implies must also be a chord in the original $\gamma$.

**Lemma 14.23** *Let $(s_1, \ldots, s_n, s_{n+1} = s_1)$ be a loop in a triangulated graph, with $n \geq 4$. Then the path has a chord at two non-contiguous vertexes at least.*

To prove the lemma, assume the contrary and let $(s_1, \ldots, s_n, s_{n+1} = s_1)$ be a loop that does not satisfy the condition, with $n$ as small as possible. If $n > 4$, the loop must have a chord, say at $s_j$, and one can remove $s_j$ from the loop to still obtain a smaller loop that must satisfy the condition in the lemma, since $n$ was as small as possible. One of the two chords must be at a vertex other than the two neighbors of $s_j$, and thus provide a second chord in the original loop, which is a contradiction. Thus $n = 4$, but $G$ being triangulated implies that this 4-point loop has a diagonal, so that the condition in the lemma also holds, which provides a contradiction.

For the "only if" part of proposition 14.22, assume that $G$ is triangulated. We prove that the graph is decomposable by induction on $|G|$. The induction will work if we can show that, if $G$ is triangulated, it is either complete or there exists a clique in $G$ such that $V \setminus C$ is disconnected, i.e., there exist two elements $a, b \in V \setminus C$ which are related by no path in $V \setminus C$. Indeed, we will then be able to decompose $V = A \cup B \cup C$, where $A$ and $B$ are unions of (distinct) connected components of $V \setminus C$. Take, for example, $A$ to be the set of vertexes connected to $A$ in $G \setminus C$, and $B = V \setminus (A \cup C)$, which is not empty since it contains $b$. Note that restricted graphs from triangulated graphs are triangulated too.

So, assume that $G$ is triangulated, and not complete. Let $C$ be a subset of $V$ that satisfies the property that $V \setminus C$ is disconnected, and take $C$ minimal, so that $V \setminus C'$ is connected for any $C' \subset C$, $C' \neq C$. We want to show that $C$ is a clique, so take $s$ and $t$ in $C$ and assume that they are not neighbors to reach a contradiction.

Let $A$ and $B$ be two connected components of $V \setminus C$. For any $a \in A$, $b \in B$, and $s, t \in C$, we know that there exists a path between $a$ and $b$ in $V \setminus C \cup \{s\}$ and another one in $V \setminus C \cup \{t\}$, the first one passing by $s$ (because it would otherwise connect $a$ and $b$ in $V \setminus C$) and the second one passing by $t$. Any point before $s$ (or $t$) in these paths must belong to $A$, and any point after them must belong to $B$. Concatenating these two paths, and removing multiple points if needed, we obtain a loop passing in $A$, then by $s$, then in $B$, then by $t$. We can recursively remove all points at which these paths have a chord. We can also notice that we cannot remove $s$ nor $t$ in this process, since this would imply an edge between $A$ and $B$, and that we must leave at least one element in $A$ and one in $B$ because removing the last one would require $s \sim t$. So, at the end, we obtain an achordal loop with at least four points, which contradicts the fact that $G$ is triangulated. ∎

We can now characterize graphs that admit junction trees over the set of their maximal cliques.

**Theorem 14.24** *Let $G = (V, E)$ be an undirected graph, and $\mathcal{C}_G^*$ be the set of all maximum cliques in $G$. The following two properties are equivalent.*

*(i) There exists a junction tree over $\mathcal{C}_G^*$.*

*(ii) G is triangulated/decomposable.*

PROOF The proof works by induction on the number of maximal cliques, $|\mathcal{C}_G^*|$. If $G$ has only one maximal clique, then $G$ is complete, because any point not included in this clique will have to be included in another maximal clique, which leads to a contradiction. So $G$ is decomposable, and, since any single node obviously provides a junction tree, (i) is true also.

Now, fix $G$ and assume that the theorem is true for any graph with fewer maximal cliques. First assume that $\mathcal{C}_G^*$ has a junction tree, $\mathbb{T}$. Let $C_1$ be a leaf in $\mathbb{T}$, connected, say, to $C_2$, and let $\mathbb{T}_2$ be $\mathbb{T}$ restricted to $\mathcal{C}_2 = \mathcal{C}_G^* \setminus \{C_1\}$. Let $V_2$ be the unions of maximal cliques from nodes in $\mathbb{T}_2$. A maximal clique $C$ in $G_{V_2}$ is a clique in $G_V$ and therefore included in some maximal clique $C' \in \mathcal{C}_V$. If $C' \in \mathcal{C}_2$, then $C'$ is also a clique in $G_{V_2}$, and for $C$ to be maximal, we need $C = C'$. If $C' = C_1$, we note that we must also have

$$C = \bigcup_{\tilde{C} \in \mathcal{C}_2} C \cap \tilde{C}$$

and whenever $C \cap \tilde{C}$ is not empty, this set must be included in any node in the path in $\mathbb{T}$ that links $\tilde{C}$ to $C_1$. Since this path contains $C_2$, we have $C \cap \tilde{C} \subset C_2$ so that $C \subset C_2$, but, since $C$ is maximal, this would imply that $C = C_2 = C_1$ which is impossible.

This shows that $\mathcal{C}_{G_2}^* = \mathcal{C}_2$. This also shows that $\mathbb{T}_2$ is a junction tree over $\mathcal{C}_2$. So, by the induction hypothesis, $G_{V_2}$ is decomposable. If $s \in V_2 \cap C_1$, then $s$ also belongs to some clique $C' \in \mathcal{C}_2$, and therefore belongs to any clique in the path between $C'$ and $C_1$, which includes $C_2$. So $s \in C_1 \cap C_2$ and $C_1 \cap V_2 = C_1 \cap C_2$. So, letting $A = C_1 \setminus (C_1 \cap C_2)$, $B = V_1 \setminus (C_1 \cap C_2)$, $S = C_1 \cap C_2$, we know that $G_{A \cup S}$ and $G_{B \cup S}$ are decomposable (the first one being complete), and that $S$ is a clique. To show that $G$ is decomposable, it remains to show that $S$ separates $A$ from $B$.

If a path connects $A$ to $B$ in $G$, it must contain an edge, say $\{s, t\}$, with $s \in V \setminus S$ and $t \in S$; $\{s, t\}$ must be included in a maximal clique in $G$. If this clique is $C_1$, we have $s \in C_1 \cap V_2 = S$. The same argument shows that this is the only possibility, because, if $\{s, t\}$ is included in some maximal clique in $\mathcal{C}_2$, then we would find $t \in C_1 \cap C_2$. So $S$ separates $A$ and $B$ in $G$.

Let us now prove the converse statement, and assume that $G$ is decomposable. If $G$ is complete, it has only one maximal clique and we are done. Otherwise, there exists a partition $V = A \cup B \cup S$ such that $G_{A \cup S}$ and $G_{B \cup S}$ are decomposable, $A$ and $B$ separated by $S$ which is complete. Let $\mathcal{C}_A^*$ be the maximal cliques in $G_{A \cup S}$ and $\mathcal{C}_B^*$ the maximal cliques in $G_{B \cup S}$. By hypothesis, there exist junction trees $\mathbb{T}_A$ and $\mathbb{T}_B$ over $\mathcal{C}_A^*$ and $\mathcal{C}_B^*$.

Let $C$ be a maximal clique in $G_{A \cup S}$. Assume that $C$ intersect $A$; $C$ can be extended to a maximal clique, $C'$, in $G$, but $C'$ cannot intersect $B$ (since this would imply a

direct edge between $A$ and $B$) and is therefore included in $A \cup S$, so that $C = C'$. Similarly, all maximal cliques in $G_{B \cup S}$ that intersect $B$ also are maximal cliques in $G$.

The clique $S$ is included in some maximal clique $S_A^* \in \mathcal{C}_A^*$. From the previous discussion, we have either $S_A^* = S$ or $S_A^* \in \mathcal{C}_G^*$. Similarly, $S$ can be extended to a maximal clique $S_B^* \in \mathcal{C}_B^*$, with $S_B^* = S$ or $S_B^* \in \mathcal{C}_G^*$. Notice also that at least one of $S_A^*$ or $S_B^*$ must be a maximal clique in $G$: indeed, assume that both sets are equal to $S$, which, as a clique, can extended to a maximal clique $S^*$ in $G$; $S^*$ must be included either in $A \cup S$ or in $B \cup S$, and therefore be a maximal clique in the corresponding graph which yields $S^* = S$. Reversing the notation if needed, we will assume that $S_A^* \in \mathcal{C}_G^*$.

All elements of $\mathcal{C}_G^*$ must belong either to $\mathcal{C}_A^*$ or $\mathcal{C}_B^*$ since any maximal clique, say $C$, in $G$ must be included in either $A \cup S$ or $B \cup S$, and therefore also provide a maximal clique in the related graph. So the nodes in $\mathbb{T}_A$ and $\mathbb{T}_B$ enumerate all maximal cliques in $G$, and we can build a tree $\mathbb{T}$ over $\mathcal{C}_G^*$ by identifying $S_A^*$ and $S_B^*$ to $S^*$ and merging the two trees at this node. To conclude our proof, it only remains to show that the running intersection property is satisfied. So consider two nodes $C, C'$ in $\mathbb{T}$ and take $s \in C \cap C'$. If the path between these nodes remain in $\mathcal{C}_A^*$, or in $\mathcal{C}_B^*$, then $s$ will belong to any set along that path, since the running intersection is true on $\mathbb{T}_A$ and $\mathbb{T}_B$. Otherwise, we must have $s \in S$, and the path must contain $S^*$ to switch trees, and $s$ must still belong to any clique in the path (applying the running intersection property between the beginning of the path and $S^*$, and between $S^*$ and the end of the path). ∎

This theorem delineates a strategy in order to build a junction tree that is adapted to a given family of local interactions $\Phi = (\varphi_C, C \in \mathcal{C})$. Letting $G$ be the graph induced by these interactions, i.e., $s \sim_G t$ if and only if there exists $C \in \mathcal{C}$ such that $\{s, t\} \subset C$, the method proceeds as follows.

(JT1) Extend $G$ by adding edges to obtain a triangulated graph $G^*$.

(JT2) Compute the set $\mathcal{C}^*$ of maximal cliques in $G^*$, which therefore extend $\mathcal{C}$.

(JT3) Build a junction tree over $\mathcal{C}^*$.

(JT4) Assign interaction $\varphi_C$ to a clique $C^* \in \mathcal{C}^*$ such that $C \subset C^*$.

(JT5) Run the junction-tree belief propagation algorithm to compute the marginal of $\pi$ (associated to $\Phi$) over each set $C^* \in \mathcal{C}^*$.

Steps (JT4) and (JT5) have already been discussed, and we now explain how the first three steps can be implemented.

### 14.6.2 Building triangulated graphs

First consider step (JT1). To triangulate a graph $G = (V, E)$, it suffices to order its vertexes so that $V = \{s_1, \ldots, s_n\}$, and then run the following algorithm.

---

**Algorithm 14.6 (Graph triangulation)**
Initialize the algorithm with $k = n$ and $E_k = E$. Given $E_k$, determine $E_{k-1}$ as follows:

- Add an edge to any pair of neighbors of $s_k$ (unless, of course, they are already linked).
- Let $E_{k-1}$ be the new set of edges.

---

Then the graph $G^* = (V, E_0)$ is triangulated. Indeed taking any achordal loop, and selecting the vertex with highest index in the loop, say $s_k$, brings a contradiction, since the neighbors of $s_k$ have been linked when building $E_{k-1}$.

However, the quality of the triangulation, which can be measured by the number of added edges, or by the size of the maximal cliques, highly depends on the way vertexes have been numbered. Take the simple example of the linear graph with three vertexes $A \sim B \sim C$. If the point of highest index is $B$, then the previous algorithm will return the three-point loop $A \sim B \sim C \sim A$. Any other ordering will leave the linear graph, which is already triangulated, invariant.

So, one must be careful about the order with which nodes will be processed. Finding an optimal ordering for a given global cost is an NP-complete problem. However, a very simple modification of the previous algorithm, which starts with $s_n$ having the minimal number of neighbors, and at each step defines $s_k$ to be the one with fewest neighbors that haven't been visited yet, provides an efficient way for building triangulations. (It has the merit of leaving $G$ invariant if it is a tree, for example). Another criterion may be preferred to the number of neighbors (for example, the number of new edges that would be needed if $s$ is added).

If $G$ is triangulated, there exists an ordering of $V$ such that the algorithm above leaves $G$ invariant. We now proceed to a proof of this statement and also show that such an ordering can be computed using an algorithm called maximum cardinality search, which, in addition, allows one to decide whether a graph is triangulated. We start with a definition that formalizes the sequence of operations in the triangulation algorithm.

**Definition 14.25** *Let $G = (V, E)$ be an undirected graph. A node elimination consists in selecting a vertex $s \in V$ and building the graph $G^{(s)} = (V^{(s)}, E^{(s)})$ with $V^{(s)} = V \setminus \{s\}$, and $E^{(s)}$ containing all pairs $\{t, t'\} \subset V^{(s)}$ such that either $\{t, t'\} \in E$ or $\{t, t'\} \subset \mathcal{V}_s$.*

$G^{(s)}$ *is called the s-elimination graph of G. The set of added edges, namely* $E^{(s)} \backslash (E \cap E^{(s)})$
*is called the deficiency set of s and denoted* $D(s)$ *(or* $D_G(s)$*).*

So, the triangulation algorithm implements a sequence of node eliminations, successively applied to $s_n, s_{n-1}$, etc. One says that such an elimination process is *perfect* if, for all $k = 1, \ldots, n$, the deficiency set of $s_k$ in the graph obtained after elimination of $s_n, \ldots, s_{k+1}$ is empty (so that no edge is added during the process). We will also say that $(s_1, \ldots, s_n)$ provides a perfect ordering for $G$.

**Theorem 14.26** *An undirected graph* $G = (V, E)$ *admits a perfect ordering if and only if it is triangulated.*

PROOF The "only if" part is obvious, since, the triangulation algorithm following a perfect ordering does not add any edge to $G$, which must therefore have been triangulated to start with.

We now proceed to the "if" part. For this it suffices to prove that for any triangulated graph, there exists a vertex $s$ such that $D_G(s) = \emptyset$. One can then easily prove the result by induction, since, after removing this $s$, the remaining graph $G^{(s)}$ is still triangulated and would admit (by induction) a perfect ordering that completes this first step.

To prove that such an $s$ exists, we take a decomposition $V = A \cup S \cup B$, in which $S$ is complete and separates $A$ and $B$, such that $|A \cup S|$ is minimal (or $|B|$ maximal). We claim that $A \cup S$ must be complete. Otherwise, since $A \cup S$ is still triangulated, There exists a similar decomposition $A \cup S = A' \cup S' \cup B'$. One cannot have $S \cap A'$ and $S \cap B'$ non empty simultaneously, since this would imply a direct edge from $A'$ to $B'$ ($S$ is complete). Say that $S \cap A' = \emptyset$, so that $A' \subset A$. Then the decomposition $V = A' \cup S' \cup (B' \cup B)$ is such that $S'$ separates $A'$ from $B \cup B'$. Indeed, a path from $A'$ to $b \in B \cup B'$ must pass in $S'$ if $b \in B'$, and, if $b \in B$, it must pass in $S$ (since it links $A$ and $B$). But $S \subset S' \cup B'$ so that the path must intersect $S'$. We therefore obtain a decomposition that enlarges $B$, which is a contradiction and shows that $A \cup S$ is complete. Given this, any element $s \in A$ can only have neighbors in $A \cup S$ and is therefore such that $D_G(s) = \emptyset$, which concludes the proof. ∎

If a graph is triangulated, there is more than one perfect ordering of its vertexes. One of these orderings is provided the *maximum cardinality search* algorithm, which also allows one to decide whether the graph is triangulated. We start with a definition/notation.

**Definition 14.27** *If* $G = (V, E)$ *is an undirected graph, with* $|V| = n$*, any ordering* $V = (s_1, \ldots, s_n)$ *can be identified with the bijection* $\alpha : V \to \{1, \ldots, n\}$ *defined by* $\alpha(s_k) = k$*. In other terms,* $\alpha(s)$ *is the rank of s in the ordering. We will refer to* $\alpha$ *as an ordering, too.*

*Given an ordering $\alpha$, we define incremental neighborhoods $\mathcal{V}_s^{\alpha,k}$, for $s \in V$ and $k = 1,\ldots,n$ to be the intersections of $\mathcal{V}_s$ with the sets $\alpha^{-1}(\{1,\ldots,k\})$, i.e.,*

$$\mathcal{V}_s^{\alpha,k} = \{t \in V, t \sim s, \alpha(t) \leq k\}.$$

*One says that $\alpha$ satisfies the maximum cardinality property if, for all $k = \{2,\ldots,n\}$*

$$|\mathcal{V}_{s_k}^{\alpha,k-1}| = \max_{\alpha(s)\geq k} |\mathcal{V}_s^{\alpha,k-1}|. \tag{14.35}$$

*where $s_k = \alpha^{-1}(k)$.*

Given this, we have the proposition:

**Proposition 14.28** *If $G = (V,E)$ is triangulated, then any ordering that satisfies the maximum cardinality property is perfect.*

Equation (14.35) immediately provides an algorithm that constructs an ordering satisfying the maximum cardinality property given a graph $G$. From proposition 14.28, we see that, if for some $k$, the largest set $\mathcal{V}_{s_k}^{\alpha,k-1}$ is not a clique, then $G$ is not triangulated. We now proceed to the proof of this proposition.

PROOF Let $G$ be triangulated, and assume that $\alpha$ is an ordering that satisfies (14.35). Assume that $\alpha$ is not proper in order to reach a contradiction.

Let $k$ be the first index for which $\mathcal{V}_{s_k}^{\alpha,k-1}$ is not a clique, so that $s_k$ has two neighbors, say $t$ and $u$, such that $\alpha(t) < k$, $\alpha(u) < k$ and $t \nsim u$. Assume that $\alpha(t) > \alpha(u)$. Then $t$ must have a neighbor that is not neighbor of $s$, say $t'$, such that $\alpha(t') < \alpha(t)$ (otherwise, $s$ would have more neighbors than $t$ at order less than $\alpha(t)$, which contradicts the maximum cardinality property). The sequence $t',t,s,u$ forms a path that is such that $\alpha$ increases from $t'$ to $s$, then decreases from $s$ to $u$, and contains no chord. Moreover, $t'$ and $u$ cannot be neighbors, since this would yield an achordal loop and a contradiction. The proof of proposition 14.28 consists in showing that this construction can be iterated until a contradiction is reached.

More precisely, assume that an achordal path $s_1,\ldots,s_k$ has been obtained, such that $\alpha(s)$ is first increasing, then decreasing along the path, and such that, at extremities one either has $\alpha(s_1) < \alpha(s_k) < \alpha(s_2)$ or $\alpha(s_k) < \alpha(s_1) < \alpha(s_{k-1})$. In fact, one can switch between these last two cases by reordering the path backwards. Both paths $(u,s,t)$ and $(u,s,t,t')$ in the discussion above satisfy this property.

• Assume, without loss of generality, that $\alpha(s_1) < \alpha(s_k) < \alpha(s_2)$ and note that, in the considered path, $s_1$ and $s_k$ cannot be neighbors (for, if $j$ is the last index smaller than $k-1$ such that $s_j$ and $s_k$ are neighbors, then $j$ must also be smaller than $k-2$ and the loop $s_j,\ldots,s_{k-1},s_k$ would be achordal).

• Since $\alpha(s_2) > \alpha(s_k)$, and $s_1$ and $s_2$ are neighbors, $s_k$ must have a neighbor, say $s_k'$, such that $s_k'$ is not neighbor of $s_2$ and $\alpha(s_k') < \alpha(s_k)$.

• Select the first index $j > 2$ such that $s_j \sim s_k'$, and consider the path $(s_1, \ldots, s_j, s_k')$. This path is achordal, by construction, and one cannot have $s_1 \sim s_k'$ since this would create an achordal loop. Let us show that $\alpha$ first increases and then decreases along this path. Since $s_2$ is in the path, $\alpha$ must first increase, and it suffices to show that $\alpha(s_k') < \alpha(s_j)$. If $\alpha$ increases from $s_1$ to $s_j$, then $\alpha(s_j) > \alpha(s_2) > \alpha(s_k) > \alpha(s_k')$. If $\alpha$ started decreasing at some point before $s_j$, then $\alpha(s_j) > \alpha(s_k) > \alpha(s_k')$.

• Finally, we need to show that the $\alpha$-value at one extremity is between the first two $\alpha$-values on the other end of the path. If $\alpha(s_k') < \alpha(s_1)$, and since we have just seen that $\alpha(s_j) > \alpha(s_k) > \alpha(s_1)$, we do get $\alpha(s_k') < \alpha(s_1) < \alpha(s_j)$. If $\alpha(s_k') > \alpha(s_1)$, then, since by construction $\alpha(s_2) > \alpha(s_k) > \alpha(s_k')$, we have $\alpha(s_2) > \alpha(s_k') > \alpha(s_1)$.

• So, we have obtained a new path that satisfies the same property that the one we started with, but with a maximum value at end points smaller than the initial one, i.e.,

$$\max(\alpha(s_1), \alpha(s_k')) < \max(\alpha(s_1), \alpha(s_k)).$$

Since $\alpha$ takes a finite number of values, this process cannot be iterated indefinitely, which yields our contradiction. ∎

### 14.6.3   Computing maximal cliques

At this point, we know that a graph must be triangulated for its maximal cliques to admit junction trees, and we have an algorithm to decide whether a graph is triangulated, and extend it into a triangulated one if needed. This provides the first step, (JT1), of our description of the junction tree algorithm. The next step, (JT2), requires computing a list of maximal cliques. Computing maximal cliques in general graph is an NP complete problem, for which a large number of algorithms has been developed (see, for example, [149] for a review). For graphs with a perfect ordering, however, this problem can always be solved in a polynomial time.

Indeed, assume that a perfect ordering is given for $G = (V, E)$, so that $V = \{s_1, \ldots, s_n\}$ is such that, for all $k$, $\mathcal{V}_{s_k}' := \mathcal{V}_{s_k} \cap \{s_1, \ldots, s_{k-1}\}$ is a clique. Let $G_k$ be $G$ restricted to $\{s_1, \ldots, s_k\}$ and $\mathcal{C}_k^*$ be the set of maximal cliques in $G_k$. Then the set $C_k := \{s_k\} \cup \mathcal{V}_{s_k}'$ is the only maximal clique in $G_k$ that contains $s_k$: it is a clique because the ordering is perfect, and any clique that contains $s_k$ must be included in it (because its elements are either $s_k$ or neighbors of $s_k$). It follows from this that the set $\mathcal{C}_k^*$ can be deduced from $\mathcal{C}_{k-1}^*$ by

$$\begin{cases} \mathcal{C}_k^* = \mathcal{C}_{k-1}^* \cup \{C_k\} \text{ if } \mathcal{V}_k' \notin \mathcal{C}_{k-1}^* \\ \mathcal{C}_k^* = (\mathcal{C}_{k-1}^* \cup \{C_k\}) \setminus \{\mathcal{V}_k'\} \text{ if } \mathcal{V}_k' \in \mathcal{C}_{k-1}^* \end{cases}$$

This allows one to enumerate all elements in $\mathcal{C}_G^* = \mathcal{C}_n^*$, starting with $\mathcal{C}_1^* = \{\{s_1\}\}$.

### 14.6.4 Characterization of junction trees

We now discuss the last remaining point, (JT3). For this, we need to form the *clique graph* of $G$, which is the undirected graph $\mathbb{G} = (\mathcal{C}_G^*, \mathbb{E})$ defined by $(C, C') \in \mathbb{E}$ if and only if $C \cap C' \neq \emptyset$. We then have the following fact:

**Proposition 14.29** *The clique graph $\mathbb{G}$ of a connected triangulated undirected graph $G$ is connected.*

Proof We proceed by induction, and assume that the result is true if $|V| = n - 1$ (the proposition obviously holds if $|V| = 1$). Assume that a perfect order on $G$ has been chosen, say $V = \{s_1, \ldots, s_n\}$. Let $G'$ be $G$ restricted to $\{s_1, \ldots, s_{n-1}\}$, and $\mathbb{G}'$ the associated clique graph. Because $\{s_n\} \cup \mathcal{V}_{s_n}$ is a clique, any path in $G$ provides a valid path in $G'$ after removing all occurrences of $s_n$ (because any two neighbors of $s_n$ are linked). The induction hypothesis also implies that $\mathbb{G}'$ is connected. Since $G$ is connected, $\mathcal{V}_{s_n}$ is not empty. Moreover, $C := \{s_n\} \cup \mathcal{V}_{s_n}$ must be a maximal clique in $G$ (since we assume that the order is perfect) and it is the only maximal clique in $G$ that contains $s_n$ (all other maximal cliques in $G$ therefore are maximal cliques in $G'$ also). To prove that $\mathbb{G}$ is connected, it suffices to prove that $C$ is connected to any other maximal clique, $C'$, in $G$ by a path in $\mathbb{G}$. If $t \in C$, $t \neq s_n$, there exists a maximal clique, say $C''$, in $G'$ that contains $t$, and, since $\mathbb{G}'$ is connected, there exists a path $(C_1 = C', \ldots, C_q = C'')$ connecting $C'$ to $C''$ in $\mathbb{G}'$. Let $j$ be the first integer such that $C_j = \mathcal{V}_n$ (take $j = q + 1$ if this never happens). Then $(C_1, \ldots, C_{j-1}, C)$ is a path linking $C'$ and $C$ in $\mathbb{G}$. ∎

We hereafter assume that $G$, and hence $\mathbb{G}$, is connected. This is not real loss of generality because connected components in undirected graphs yields independent processes that can be handled separately. We assign weights to edges of the clique graph of $G$ by defining $w(C, C') = |C \cap C'|$. A subgraph $\tilde{T}$ of any given graph $\tilde{G}$ is called a spanning tree if $\tilde{T}$ is a tree with set of vertexes equal to the set of vertexes of $\tilde{G}$. If $\mathbb{T} = (\mathcal{C}_G^*, \mathbb{E}')$ is a spaning tree of $\mathbb{G}$, we define the total weight

$$w(\mathbb{T}) = \sum_{\{C, C'\} \in \mathbb{E}'} w(C, C').$$

We then have the proposition:

**Proposition 14.30** *[99] If $G$ is a connected triangulated graph, the set of junction trees over $\mathcal{C}_G^*$ coincides with the set of maximizers of $w(\mathbb{T})$ over all spanning trees of $\mathbb{G}$.*

(Notice that $\mathbb{G}$ being connected implies that spanning trees over $\mathbb{G}$ exist.)

Before proving this proposition, we discuss some properties related to maximal (or maximum-weight) spanning trees over an undirected graph. For this discussion, we let $G = (V, E)$ be any undirected graph with weight $(w(e), e \in E)$. We will then apply these results to a clique graph when will switch back to the general notation of this section. Maximal spanning trees can be computed using the so-called Prim's algorithm [98, 155, 63].

---

**Algorithm 14.7 (Prim's algorithm)**
Initialize the algorithm with a single-node tree $T_1 = (\{s_1\}, \emptyset)$, for some arbitrary $s_1 \in V$. Let $T_{k-1} = (V_{k-1}, E_{k-1})$ be the tree obtained at step $k - 1$ of the algorithm. If $k \leq n$, the next tree is built as follows.

(1) Let
$$V_k = \{s_k\} \cup V_{k-1} \ (s_k \notin V_{k-1}.)$$

(2) Let $E_k = \{e_k\} \cup E_{k-1}$, such that $e_k = \{s_k, s\}$ for some $s \in V_{k-1}$ satisfying

$$w(e_k) = \max\left(w(\{t, t'\}), \{t, t'\} \in E, t \notin V_{k-1}, t' \in V_{k-1}\right). \tag{14.36}$$

---

The ability of this algorithm to always build a maximal spanning tree is summarized in the following proposition [81, 129].

**Proposition 14.31** *If $G = (V, E, w)$ is a weighted, connected undirected graph, Prim's algorithm, as described above, provides a sequence $T_k = (V_k, E_k)$, for $k = 1, \ldots, n$ of subtrees of $G$ such that $V_n = V$ and, for all $k$, $T_k$ is a maximal spanning tree for the restriction $G_{V_k}$ of $G$ to $V_k$.*

*Moreover, any maximal spanning tree of $G$, can be realized as $T_n$, where $(T_1, \ldots, T_n)$ is a sequence provided by Prim's algorithm.*

PROOF We first prove that, for all $k$, $T_k$ is a maximal spanning tree on the graph $G_{V_k}$.

We will prove a slightly stronger statement, namely, that, for all $k$, $T_k$ can be extended to form a maximal spanning tree of $G$. This is stronger, because, if $T_k = (V_k, E_k)$ can be extended to a maximal spanning tree $T = (V, E)$, and if $T'_k = (V_k, E'_k)$ is a spanning tree for $G_{V_k}$ such that $w(T_k) < w(T'_k)$, then the graph $T' = (V, E')$ with

$$E' = (E \setminus E_k) \cup E'_k$$

would be a spanning tree for $G$ with $w(T) < w(T')$, which is impossible. (To see that $T'$ is a tree, notice that paths in $T'$ are in one-to-one correspondence with paths in $T$ by replacing any subpath within $T'_k$ by the unique subpath in $T_k$ that has the same extremities.)

Clearly, $T_1$, which only has one vertex, can be extended to a maximal spanning tree. Let $k \geq 1$ be the last integer for which this property is true for all $j = 1, \ldots, k$. If $k = n$, we are done. Otherwise, take a maximum spanning tree, $T$, that extends $T_k$. This tree cannot contain the new edge added when building $T_{k+1}$, namely $e_{k+1} = \{s_{k+1}, s\}$ as defined in Prim's algorithm, since it would otherwise also extend $T_{k+1}$. Consider the path $\gamma$ in $T$ that links $s$ to $s_k$. This path must have an edge $e = \{t, t'\}$ such that $t \in V_k$ and $t' \notin V_k$, and by definition of $e_{k+1}$, we must have $w(e) \leq w(e_{k+1})$. Notice that $e$ is uniquely defined, because a path leaving $V_k$ cannot return in this set, since one would be otherwise able to close it into a loop by inserting the only path in $T_k$ that connects its extremities.

Replace $e$ by $e_{k+1}$ in $T$. The resulting graph, say $T'$, is still a spanning tree for $G$. From any path in $T$, one can create a path in $T'$ with the same extremities by replacing any occurrence of the edge, $e$, by the concatenation of the unique path in $T$ going from $t$ to $s$, followed by $(s, s_{k+1})$, followed by the unique path in $T$ going from $s_{k+1}$ to $t'$. This implies that $T'$ is connected. It is also acyclic, since any loop in $T$ would have to contain $e_{k+1}$ (since $T$ is acyclic), but there is no other path than $(s, s_{k+1})$ in $T'$ that links $s$ and $s_k$, because this path would have to be in $T$, and we have removed the only possible one from $T$ by deleting the edge $e$.

As a conclusion, $T'$ is an extension of $T_{k+1}$, and a spanning tree with total weight larger or equal to the one of $T$, and must therefore be optimal, too. But this contradicts the fact that $T_{k+1}$ cannot be extended to a maximal tree, so that $k = n$ and the sequence of trees provided by Prim's algorithm is optimal.

To prove the second statement, let $T$ be an optimal spanning tree. Let $k$ be the largest integer such that there exists a sequence $(T_1, \ldots, T_k)$ generated by Prim's algorithm, such that, for all $j = 1, \ldots, k$, $T_j$ is a subtree of $T$. One necessarily has $j \geq 1$, since $T$ extends any one-vertex tree. If $k = n$, we are done. Assuming otherwise, let $T_k = (V_k, E_k)$ and make one more step of Prim's algorithm, selecting an edge $e_{k+1} = (s_{k+1}, s)$ satisfying (14.36). By assumption, $e_{k+1}$ is not in $T$. Take as before the unique path linking $s$ and $s_{k+1}$ in $T$ and let $e$ be the unique edge at which this path leaves $V_k$. Replacing $e$ by $e_{k+1}$ in $T$ provides a new spanning tree, $T'$. One must have $w(e) \geq w(e_{k+1})$ because $T$ is optimal, and $w(e_{k+1}) \geq w(e)$ by (14.36). So $w(e) = w(e_{k+1})$, and one can use $e$ instead of $e_{k+1}$ for the $(k+1)$th step of Prim's algorithm. But this contradicts the fact that $k$ was the largest integer in a sequence of subtrees of $T$ that is generated by Prim's algorithm, and one therefore has $k = n$. ■

The proof of proposition 14.30, that we provide now, uses very similar "edge-switching" arguments.

Proof (Proof of proposition 14.30) Let us start with a maximum weight spanning tree for $\mathbb{G}$, say $\mathbb{T}$, and show that it is a junction tree. Since $\mathbb{T}$ has maximum weight, we

know that it can be obtained via Prim's algorithm, and that there exists a sequence $\mathbb{T}_1, \ldots, \mathbb{T}_n = \mathbb{T}$ of trees constructed by this algorithm. Let $\mathbb{T}_k = (\mathcal{C}_k, \mathbb{E}_k)$.

We proceed by contradiction. Let $k$ be the largest index such that $\mathbb{T}_k$ can be extended to a junction tree for $\mathcal{C}_G^*$, and let $\mathbb{T}'$ be a junction tree extension of $\mathbb{T}_k$. Assume that $k < n$, and let $e_{k+1} = (C_{k+1}, C')$ be the edge that has been added when building $\mathbb{T}_{k+1}$, with $\mathcal{C}_{k+1} = \{C_{k+1}\} \cup \mathcal{C}_k$. This edge is not in $\mathbb{T}'$, so that there exists a unique edge $e = (B, B')$ in the path between $C_k$ and $C'$ in $\mathbb{T}'$ such that $B \in \mathcal{C}_k$ and $B' \notin \mathcal{C}_k$. We must have $w(e) = |B \cap B'| \leq w(e_{k+1}) = |C_{k+1} \cap C'|$. But, since the running intersection property is true for $\mathbb{T}'$, both $B$ and $B'$ must contain $C_{k+1} \cap C'$ so that $B \cap B' = C_{k+1} \cap C'$. This implies that, if one modifies $\mathbb{T}'$ by replacing edge $e$ by edge $e_{k+1}$, yielding a new spanning tree $\mathbb{T}''$, the running intersection property is still satisfied in $\mathbb{T}'$. Indeed if a vertex $s \in V$ belongs to both extremities of a path containing $B$ and $B'$ in $\mathbb{T}'$, then it must belong to $B \cap B'$, and hence to $C_{k+1} \cap C'$, and therefore to any set in the path in $\mathbb{T}'$ that linked $C_{k+1}$ and $C'$. So we found a junction tree extension of $\mathbb{T}_{k+1}$, which contradicts our assumption that $k$ was the largest. We must therefore have $k = n$ and $\mathbb{T}$ is a junction tree.

Let us now consider the converse statement and assume that $\mathbb{T}$ is a junction tree. Let $k$ be the largest integer such that there exists a sequence of subgraphs of $\mathbb{T}$ that is provided by Prim's algorithm. Denote such a sequence by $(\mathbb{T}_1, \ldots, \mathbb{T}_k)$, with $\mathbb{T}_j = (\mathcal{C}_j, \mathbb{E}_j)$. Assume (to get a contradiction) that $k < n$, and consider a new step for Prim's algorithm, adding a new edge $e_{k+1} = \{C_{k+1}, C'\}$ to $\mathbb{T}_k$. Take as before the path in $\mathbb{T}$ linking $C'$ to $C_{k+1}$ in $\mathbb{T}$, and select the edge $e$ at which this path leaves $\mathcal{C}_k$. If $e = (B, B')$, we must have $w(e) = |B \cap B'| \leq w(e_k) = |C_{k+1} \cap C'|$, and the running intersection property in $\mathbb{T}$ implies that $C_{k+1} \cap C' \subset B \cap B'$, which implies that $w(e) = w(e_{k+1})$. This implies that adding $e$ instead of $e_{k+1}$ at step $k+1$ is a valid choice for Prim's algorithm, and contradicts the fact that $k$ was the largest number of such steps that could provide a subtree of $\mathbb{T}$. So $k = n$ and $\mathbb{T}$ is maximal.      ∎

# Chapter 15

# Bayesian Networks

## 15.1 Definitions

Bayesian networks are graphical models supported by directed acyclic graphs (DAG), which provide them with an ordered organization (directed graphs were introduced in definition 13.35).

We first introduce some notation. Let $G = (V, E)$ be a directed acyclic graph. The parents of $s \in V$ are vertexes $t$ such that $(t, s) \in E$, and its children are $t$'s such that $(s, t) \in E$. The set of parents of $s$ is denoted $pa(s)$, and the set of its children is $ch(s)$, with $\mathcal{V}_s = ch(s) \cup pa(s)$.

Similarly to trees, the vertexes of $G$ can be partially ordered by $s \leq_G t$ if and only if there exists a path going from $s$ to $t$. Unlike trees, however, there can be more than one minimal element in $V$, and we still call roots vertexes that have no parent, denoting

$$V_0 = \{s \in V : pa(s) = \emptyset\}.$$

We also call leaves, or terminal nodes, vertexes that have no children. Unless otherwise specified, we assume that all graphs are connected.

Bayesian networks over $G$ are defined as follows. We use the same notation as with Markov random fields to represent the set of configurations $\mathcal{F}(V)$ that contains collections $x = (x_s, s \in V)$ with $x_s \in F_s$.

**Definition 15.1** *A random variable X with values in $\mathcal{F}(V)$ is a Bayesian network over a DAG $G = (V, E)$ if and only if its distribution can be written in the form*

$$P_X(x) = \prod_{s \in V_0} p_s(x^{(s)}) \prod_{s \in V \setminus V_0} p_s(x^{(pa(s))}, x^{(s)}) \tag{15.1}$$

*where $p_s$ is, for all $s \in V$, a probability distribution with respect to $x^{(s)}$.*

371

Using the convention that conditional distributions given the empty set are just absolute distributions, we can rewrite (15.1) as

$$P_X(x) = \prod_{s \in V} p_s(x^{(pa(s))}, x^{(s)}). \tag{15.2}$$

One can verify that $\sum_{x \in \Omega} P^X(x) = 1$. Indeed, when summing over $x$, we can start summing over all $x^{(s)}$ with $ch(s) = \emptyset$ (the leaves). Such $x^{(s)}$'s only appear in the corresponding $p_s$'s, which disappear since they sum to 1. What remains is the sum of the product over $V$ minus the leaves, and the argument can be iterated until the remaining sum is 1 (alternatively, work by induction on $|V|$). This fact is also a consequence of proposition 15.5 below, applied with $A = \emptyset$.

## 15.2   Conditional independence graph

### 15.2.1   Moral graph

Bayesian networks have a conditional independence structure which is not exactly given by $G$, but can be deduced from it.  Indeed, fixing $S \subset V$, we can see, when computing the probability of $X^{(S)} = x^{(s)}$ given $X^{(S^c)} = x^{(S^c)}$, which is

$$\mathbb{P}(X^{(S)} = x^{(S)} \mid X^{(S^c)} = x^{(S^c)}) = \frac{1}{Z(x^{(S^c)})} \prod_{s \in V} p_s(x^{(pa(s))}, x^{(s)}),$$

that the only variables $x^{(t)}, t \notin S$ that can be factorized in the normalizing constant are those that are neither parent nor children of vertexes in $S$, and do not share a child with a vertex in $S$ (i.e., they intervene in no $p_s(x^{(pa(s))}, x^{(s)})$ that involve elements of $S$). This suggests the following definition.

**Definition 15.2** *Let $G$ be a directed acyclic graph. We denote $G^\sharp = (V, E^\sharp)$ the undirected graph on $V$ such that $\{s, t\} \in E^\sharp$ if one of the following conditions is satisfied*

- *Either $(s, t) \in E$ or $(t, s) \in E$.*

- *There exists $u \in V$ such that $(s, u) \in E$ and $(t, u) \in E$.*


$G^\sharp$ is sometimes called the *moral graph* of $G$ (because it forces parents to marry !). A path in $G^\sharp$ can be visualized as a path in $G^\flat$ (the undirected graph associated with $G$) which is allowed to jump between parents of the same vertex even if they were not connected originally.

The previous discussion implies:

**Proposition 15.3** *Let X be a Bayesian network on G. We have*

$$(S \perp\!\!\!\perp T \mid U)_{G^\sharp} \Rightarrow (X^{(S)} \perp\!\!\!\perp X^{(T)} \mid X^{(U)}),$$

*i.e., X is $G^\sharp$-Markov.*

This proposition can be refined by noticing that the joint distribution of $X^{(S)}$, $X^{(T)}$ and $X^{(U)}$ can be deduced from a Bayesian network on a graph restricted to the ancestors of $S \cup T \cup U$. Definition 13.21 for restricted graphs extends without change to directed graphs, and we repeat it below for convenience.

**Definition 15.4** *Let $G = (V,E)$ be a graph (directed or undirected), and $A \subset V$. The restricted graph $G_A = (A, E_A)$ is such that the elements of $E_A$ are the edges $(s,t)$ (or $\{s,t\}$) in E such that both s and t belong to A.*

Moreover, for a directed acyclic graph $G$ and $s \in V$, we define the set of ancestors of $s$ by

$$\mathcal{A}_s = \{t \in V, t \leq_G s\} \tag{15.3}$$

for the partial order on $V$ induced by $G$.

If $S \subset V$, we denote $\mathcal{A}_S = \bigcup_{s \in S} \mathcal{A}_s$. Note that, by definition, $S \subset \mathcal{A}_S$. The following proposition is true.

**Proposition 15.5** *Let X be a Bayesian network on $G = (V,E)$ with distribution given by (15.2). Let $S \subset V$ and $A = \mathcal{A}_S$. Then the distribution of $X^{(A)}$ is a Bayesian network over $G_A$ given by*

$$\mathbb{P}(X^{(A)} = x^{(A)}) = \prod_{s \in A} p_s(x^{(pa(s))}, x^{(s)}). \tag{15.4}$$

There is no ambiguity in the notation $pa(s)$, since the parents of $s \in A$ are the same in $G_A$ as in $G$.

PROOF One needs to show that

$$\prod_{s \in A} p_s(x^{(pa(s))}, x^{(s)}) = \sum_{x_{A^c}} \prod_{s \in V} p_s(x^{(pa(s))}, x^{(s)}).$$

This can be done by induction on the cardinality of $V$. Assume that the result is true for graphs of size $n$, and let $|V| = n+1$ (the result is obvious for graphs of size 1).

If $A = V$, there is nothing to prove, so assume that $A^c$ is not empty. Then $A^c$ must contain a leaf in $G$, since otherwise, $A$ would contain all leaves and their ancestors which would imply that $A = V$.

If $s \in A^c$ is a leaf in $G$, one can remove the variable $x^{(s)}$ from the sum, since it only appear in $p_s$ and transition probabilities sum to one. But one can now apply the induction assumption to the restriction of $G$ to $V \setminus \{s\}$. ∎

Given proposition 15.5, proposition 15.3 can therefore be refined as follows.

**Proposition 15.6** *Let X be a Bayesian network on G. We have*

$$(S \perp\!\!\!\perp T \mid U)_{(G_{\mathcal{A}_{S \cup T \cup U}})^\sharp} \Rightarrow (X^{(S)} \perp\!\!\!\perp X^{(T)} \mid X^{(U)}).$$

Proposition 15.5 is also used in the proof of the following proposition.

**Proposition 15.7** *Let $G = (V, E)$ be a directed acyclic graph, and X be a Bayesian network over G. Then, for all $s \in S$*

$$\mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(\mathcal{A}_s \setminus \{s\})} = x^{(\mathcal{A}_s \setminus \{s\})}) = \mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(pa(s))} = x^{(s^-)}) = p_s(x^{(pa(s))}, x^{(s)}).$$

Proof By proposition 15.5, we can without loss of generality assume that $V = \mathcal{A}_s$. Then

$$
\begin{aligned}
\mathbb{P}(X^{(s)} = x^{(s)} \mid X^{(\mathcal{A}_s \setminus \{s\})} = x^{(\mathcal{A}_s \setminus \{s\})}) \quad &\propto \quad \mathbb{P}(X^{(\mathcal{A}_s)} = x^{(\mathcal{A}_s)}) \\
&= \quad p_s(x^{(pa(s))}, x^{(s)}) Z(x^{(\mathcal{A}_s \setminus \{s\})})
\end{aligned}
$$

where

$$Z(x^{(\mathcal{A}_s \setminus \{s\})}) = \prod_{t \in \mathcal{A}_s \setminus \{s\}} p_t(x^{(pa(t))}, x^{(t)})$$

disappears when the conditional probability is normalized.                    ∎

### 15.2.2   Reduction to d-separation

We now want to reformulate proposition 15.6 in terms of the unoriented graph $G^\flat$ and specific features in $G$ called v-junctions, that we now define.

**Definition 15.8** *Let $G = (V, E)$ be a directed graph. A v-junction is a triple of distinct vertexes, $(s, t, u) \in V \times V \times V$ such that $\{s, u\} \subset pa(t)$ (i.e., s and u are parents of t).*

*We will say that a path $(s_1, \dots, s_N)$ in $G^\flat$ passes at $s = s_k$ with a v-junction if $(s_{k-1}, s_k, s_{k+1})$ is a v-junction in G.*

We have the lemma:

**Lemma 15.9** *Two vertexes s and t in G are separated by a set U in $(G_{\mathcal{A}_{\{s,t\} \cup U}})^\sharp$ if and only if any path between s and t in $G^\flat$ must either*

*(1)  Pass at a vertex in U without a v-junction.*

*(2)  Pass in $V \setminus \mathcal{A}_{\{s,t\} \cup U}$ at a v-junction.*

Proof

*Step 1.* We first note that the v-junction clause is redundant in (2). It can be removed without affecting the condition. Indeed, if a path in $G^\flat$ passes in $V \setminus \mathcal{A}_{\{s,t\}\cup U}$ one can follow this path downward (i.e., following the orientation in $G$) until a v-junction is met. This has to happen before reaching the extremities of the path, since $u$ would be an ancestor of $s$ or $t$ otherwise. We can therefore work with the weaker condition (that we will denote (2)') in the rest of proof.

*Step 2.* Assume that $U$ separates $s$ and $t$ in $(G_{\mathcal{A}_{\{s,t\}\cup U}})^\sharp$. Take a path $\gamma$ between $s$ and $t$ in $G^\flat$. We need to show that the path satisfies (1) or (2)'. So assume that (2)' is false (otherwise we are done) so that $\gamma$ is included in $\mathcal{A}_{\{s,t\}\cup U}$. We can modify $\gamma$ by removing all the central nodes in v-junctions and still keep a valid path in $(G_{\mathcal{A}_{\{s,t\}\cup U}})^\sharp$ (since parents are connected in the moral graph). The remaining path must intersect $U$ by assumption, and this cannot be at a v-junction in $\gamma$ since we have removed them. So (1) is true.

*Step 3.* Conversely, assume that (1) or (2) is true for any path in $G^\flat$. Consider a path $\gamma$ in $(G_{\mathcal{A}_{\{s,t\}\cup U}})^\sharp$ between $s$ and $t$. Any edge in $\gamma$ that is not in $G^\flat$ must involve parents of a common child in $\mathcal{A}_{\{s,t\}\cup U}$. Insert this child between the parents every time this occurs, resulting in a v-junction added to $\gamma$. Since the added vertexes are still in $\mathcal{A}_{\{s,t\}\cup U}$, the new path still has no intersection with $V \setminus \mathcal{A}_{\{s,t\}\cup U}$ and must therefore satisfy (1). So there must be an intersection with $U$ without a v-junction, and since the new additions are all at v-junctions, the intersection must have been originally in $\gamma$, which therefore passes in $U$. This shows that $U$ separates $s$ and $t$ in $(G_{\mathcal{A}_{\{s,t\}\cup U}})^\sharp$.∎

Condition (2) can be further restricted to provide the notion of $d$-separation.

**Definition 15.10** *One says that two vertexes $s$ and $t$ in $G$ are $d$-separated by a set $U$ if and only if any path between $s$ and $t$ in $G^\flat$ must either*

*(D1) Pass at a vertex in $U$ without a v-junction.*

*(D2) Pass in $V \setminus \mathcal{A}_U$ with a v-junction.*

Then we have:

**Theorem 15.11** *Two vertexes $s$ and $t$ in $G$ are separated by a set $U$ in $(G_{\mathcal{A}_{\{s,t\}\cup U}})^\sharp$ if and only if they are $d$-separated by $U$.*

Proof It suffices to show that if condition ((D1) or (D2)) holds for any path between $s$ and $t$ in $G^\flat$, then so does ((1) or (2)). So take a path between $s$ and $t$: if (D1) is true

for this path, the conclusion is obvious, since (D1) and (1) are the same. So assume
that (D1) (and therefore (1)) is false and that (D2) is true. Let $u$ be a vertex in $V \setminus \mathcal{A}_U$
at which $\gamma$ passes with a v-junction.

Assume that (2) is false. Then $u$ must be an ancestor of either $s$ or $t$. Say it is an
ancestor of $s$: there is a path in $G$ going from $u$ to $s$ without passing by $U$ (otherwise
$u$ would be an ancestor of $U$); one can replace the portion of the old path between
$s$ and $u$ by this new one, which does not pass by $u$ with a v-junction anymore. So
the new path still does not satisfy (D1) and must satisfy (D2). Keep on removing
all intersections with ancestors of $s$ and $t$ that have v-junctions to finally obtain a
path that satisfies neither (D1) or (D2) and a contradiction to the fact that $s$ and $t$ are
$d$-separated by $U$.                                                                      ∎

### 15.2.3   Chain-graph representation

The $d$-separability property involves both unoriented and oriented edges. It is in
fact a property of the hybrid graph in which the orientation is removed from the
edges that are not involved in a v-junction, and retained otherwise. Such graphs are
particular instances of chain graphs.

**Definition 15.12** *A chain graph $G = (V, E, \tilde{E})$ is composed with a finite set $V$ of vertexes,
a set $E \subset \mathcal{P}_2(V)$ of unoriented edges and a set $\tilde{E} \subset E \times E \setminus \{(t, t), t \in E\}$ of oriented edges
with the property that $E \cap \tilde{E}^\flat = \emptyset$, i.e., two vertexes cannot be linked by both an oriented
and an unoriented edge.*

*A path in a chain graph is a a sequence of vertexes $s_0, \ldots, s_N$ such that for all $k \geq 1$,
$s_{k-1}$ and $s_k$ form an edge, which means that either $\{s_{k-1}, s_k\} \in E$ or $(s_{k-1}, s_k) \in \tilde{E}$.*

*A chain graph is acyclic if it contains no loop. It is semi-acyclic if it contains no loop
containing oriented edges.*

We start with the following equivalence relation within vertexes in a semi-acyclic
chain graph.

**Proposition 15.13** *Let $G = (V, E, \tilde{E})$ be a semi-acyclic chain graph. Define the relation
$s \mathcal{R} t$ if and only if there exists a path in the unoriented subgraph $(V, E)$ that links $s$ and $t$.
Then $\mathcal{R}$ is an equivalence relation.*

The proposition is obvious. This relation partitions $V$ in equivalence classes, the
set of which being denoted $V_\mathcal{R}$. If $S \in V_\mathcal{R}$, then any pair $s, t$ in $S$ is related by an
unoriented path, and if $S \neq S' \in V_\mathcal{R}$, no elements $s \in S$ and $t \in S'$ can be related by
such a path.

Moreover, no path in $G$ between two elements of $S \in V_{\mathcal{R}}$, can contain a directed edge, since these elements must also be related by an undirected path, and this would create a loop in $G$ containing an undirected edge. So the restriction of $G$ to $S$ is an undirected graph.

One can define a directed graph over equivalence classes as follows. Let $G_{\mathcal{R}} = (V_{\mathcal{R}}, E_{\mathcal{R}})$ be such that $(S, S') \in E_{\mathcal{R}}$ if and only if there exists $s \in S$ and $t \in S'$ such that $(s, t) \in \tilde{E}$. The graph $G_{\mathcal{R}}$ is acyclic: any loop in $G_{\mathcal{R}}$ would induce a loop in $G$ containing at least one oriented edge.

We now can formally define a probability distribution on a semi-acyclic chain graph.

**Definition 15.14** *Let $G = (V, E, \tilde{E})$ be a semi-acyclic chain graph. One says that a random variable $X$ decomposes on $G$ if and only if: $(X^{(S)}, S \in V_{\mathcal{R}})$ is a Bayesian network on $G_{\mathcal{R}}$ and the conditional distribution of $X^{(S)}$ given $X^{(S')}, S' \in pa(S)$ is $G_S$-Markov, such that, for $s \in S$, $P(X^{(s)} = x^{(s)} \mid X^{(t)}, t \in S, X_{S'}, S' \in pa(S))$ only depends on $x^{(t)}$ with $\{s, t\} \in E$ or $(t, s) \in \tilde{E}$.*

Returning to our discussion on Bayesian networks, we have the following. Associate to a DAG $G = (V, E)$ the chain graph $G^{\dagger} = (V, E^{\dagger}, \tilde{E}^{\dagger})$ defined by: $\{s, t\} \in E^{\dagger}$ if and only if $(s, t)$ or $(t, s) \in E$ and is not involved in a v-junction, and $(s, t) \in \tilde{E}^{\dagger}$ if $(s, t) \in E$ and is involved in a v-junction. This graph is acyclic; indeed, take any loop in $G^{\dagger}$: when its edges are given their original orientations in $E$, the sequence cannot contain a v-junction, since the orientation in v-junctions are kept in $G^{\dagger}$; the path therefore constitutes a loop in $G$ which is a contradiction.

All, excepted at most one, vertexes in an equivalence class $S \in G_{\mathcal{R}}^{\dagger}$ have all their parents in $S$. Indeed, assume that two vertexes, $s$ and $t$, in $S$ have parents outside of $S$. There exists an unoriented path, $s_0 = s, s_1, \ldots, s_N = t$, in $G^{\dagger}$ connecting them, since they belong to the same equivalence class. The edge at $s$ must be oriented from $s$ to $s_1$ in $G$, since otherwise $s_1$ would be a second parent to $s$ in $G$, creating a v-junction, and the edge would have remained oriented in $G^{\dagger}$. Similarly, the last edge in the path must be oriented from $t$ to $s_{N-1}$ in $G$. But this implies that there exists a v-junction in the original orientation along the path, which cannot be constituted with only unoriented edges in $G^{\dagger}$. So we get a contradiction.

Thus, random variables that decompose on $G^{\dagger}$ are "Bayesian networks" of acyclic graphs, or trees since we know these are equivalent. The root of each tree must have multiple (vertex) parents in the parent tree in $G_{\mathcal{R}}$. The following theorem states that all Bayesian networks are equivalent to such a process.

**Theorem 15.15** *Let $G = (V, E)$ be a DAG. The random variable $X$ is a Bayesian network on $G$ if and only if it decomposes over $G^{\dagger}$.*

PROOF Assume that $X$ is a Bayesian network on $G$.  We can obviously rewrite the probability distribution of $X$ in the form

$$\pi(x) = \prod_{S \in G_{\mathcal{R}}^{\dagger}} \prod_{s \in S} p_s(x^{(pa(s))}, x^{(s)}).$$

Since every vertex in $S$ has its parents in $S$ or in $\bigcup_{T \in pa(S)} T$, this *a fortiori* takes the form

$$\pi(x) = \prod_{S \in G_{\mathcal{R}}^{\dagger}} p_S((x^{(T)}, T \in S^-), x^{(s)}).$$

So $X^{(S)}, S \in V_{\mathcal{R}}$ is a Bayesian network. Moreover,

$$p_S((x^{(T)}, T \in S^-), x^{(s)}) = \prod_{s \in S} p_s(x^{(pa(s))}, x^{(s)})$$

is a tree distribution with the required form of the individual conditional distributions.

Now assume that $X$ decomposes on $G^{\dagger}$.  Then the conditional distribution of $X^{(S)}$ given $X^{(T)}, T \in pa(S)$ is Markov for the acyclic undirected graph $G_S$, and can therefore be expressed as a tree distribution consistent with the orientation of $G$.  ∎

### 15.2.4   Markov equivalence

While the previous discussion provides a rather simple description of Bayesian networks in terms of chain graphs, it does not go all the way in reducing the number of oriented edges in the definition of a Bayesian network. The issue is, in some way, addressed by the notion of Markov equivalence, which is defined as follows.

**Definition 15.16** *Two directed acyclic graphs on the same set of vertexes $G = (V, E)$ and $\tilde{G} = (V, \tilde{E})$ are Markov-equivalent if any family of random variables that decomposes as a (positive) Bayesian network over one of them also decomposes as a Bayesian network over the other.*

The notion of Markov equivalence is exactly described by d-separation.  This is stated in the following theorem, due to Geiger and Pearl [77, 76], that we state without proof.

**Theorem 15.17** *$G$ and $\tilde{G}$ are Markov equivalent if and only if, whenever two vertexes are d-separated by a set in one of them, the same separation is true with the other.*

This property can be expressed in a strikingly simple condition. One says that a v-junction $(s, t, u)$ in a DAG is *unlinked* if $s$ and $u$ are not neighbors.

**Theorem 15.18** *G and $\tilde{G}$ are Markov equivalent if and only if $G^{\flat} = \tilde{G}^{\flat}$ and G and $\tilde{G}$ have the same unlinked v-junctions.*

PROOF *Step 1.* We first show that a given pair of vertexes in a DAG is unlinked if and only if it can be d-separated by some set in the graph. Clearly, if they are linked, they cannot be d-separated (which is the "if" part), so what really needs to be proved is that unlinked vertexes can be d-separated. Let $s$ and $t$ be these vertexes and let $U = \mathcal{A}_{\{s,t\}} \setminus \{s,t\}$. Then $U$ d-separates $s$ and $t$ since any path between $s$ and $t$ in $(G_{\mathcal{A}_{\{s,t\} \cup U}})^{\sharp} = (G_{\mathcal{A}_{\{s,t\}}})^{\sharp}$ must obviously pass in $U$.

*Step 2.* We now prove the only-if part of theorem 15.18 and therefore assume that $G$ and $\tilde{G}$ are Markov equivalent, or, as stated in theorem 15.17, that d-separation coincides in $G$ and $\tilde{G}$. We want to prove that $G^{\flat} = \tilde{G}^{\flat}$ and unlinked v-junctions are the same.

*Step 2.1.* The first statement is obvious from Step 1: d-separation determines the existence of a link, so if d-separation coincides in the two graphs, then the same holds for links and $G^{\flat} = \tilde{G}^{\flat}$.

*Step 2.2.* So let us proceed to the second statement and let $(s,t,u)$ be an unlinked v-junction in $G$. We want to show that it is also a v-junction in $\tilde{G}$ (obviously unlinked since links coincide).

We will denote by $\tilde{\mathcal{A}}_S$ the ancestors of some set $S \subset V$ in $\tilde{G}$ (while $\mathcal{A}_S$ still denotes its ancestors in $G$). Let $U = \mathcal{A}_{\{s,u\}} \setminus \{s,u\}$. Then, as we have shown in Step 1, $U$ d-separates $s$ and $u$ in $G$, so that, by assumption it also d-separates them in $\tilde{G}$.

We know that $t \notin U$, because it cannot be both a child and an ancestor of $\{s,u\}$ in $G$ (this would induce a loop). The path $(s,t,u)$ links $s$ and $u$ and does not pass in $U$, which is only possible (since $U$ d-separates $s$ and $t$ in $\tilde{G}$) if it passes in $V - \tilde{\mathcal{A}}_U$ at a v-junction: so $(s,t,u)$ is a v-junction in $\tilde{G}$, which is what we wanted to prove.

*Step 3.* We now consider the converse statement and assume that $G^{\flat} = \tilde{G}^{\flat}$ and unlinked v-junctions coincide. We want to show that d-separation is the same in $G$ and $\tilde{G}$. So, we assume that $U$ d-separates $s$ and $t$ in $G$, and we want to show that the same is true in $\tilde{G}$. Thus, what we need to prove is:

*Claim 1.* Consider a path $\gamma$ between $s$ and $t$ in $\tilde{G}^{\flat} = G^{\flat}$. Then $\gamma$ either (D1) passes in $U$ without a v-junction in $\tilde{G}$, or (D2) in $V \setminus \tilde{\mathcal{A}}_U$ with a v-junction in $\tilde{G}$.

We will prove Claim 1 using a series of lemmas. We say that $\gamma$ has a three-point loop at $u$ if $(v,u,w)$ are three consecutive points in $\gamma$ such that $v$ and $w$ are linked. So $(v,u,w,v)$ forms a loop in the undirected graph.

**Lemma 15.19** *If $\gamma$ is a path between s and t that does not satisfy (D2) for G and passes in $U$ without three-point loops, then $\gamma$ satisfies (D1) for $\tilde{G}$.*

The proof is easy: since $\gamma$ does not satisfy (D2) in $G$, it satisfies (D1) and passes in $U$ without a v-junction in $G$. But this intersection cannot be a v-junction in $\tilde{G}$ since

it would otherwise have to be linked and constitute a three-point loop in $\gamma$, which proves that (D1) is true for $\gamma$ in $\tilde{G}$.

The next step is to remove the three-point loop condition in lemma 15.19. This will be done using the next two results.

**Lemma 15.20** *Let $\gamma$ be a path with a three-point loop at $u \in U$ for G. Assume that $\gamma \setminus u$ (which is a valid path in $G^\flat$) satisfies (D1) or (D2) in $\tilde{G}$. Then $\gamma$ satisfies (D1) or (D2) in $\tilde{G}$.*

To prove the lemma, let $v$ and $w$ be the predecessor and successor of $u$ in $\gamma$. First assume that $\gamma \setminus u$ satisfies (D1) in $\tilde{G}$. If this does not happen at $v$ or at $w$, then this will apply also to $\gamma$ and we are done, so let us assume that $v \in U$ and that $(v', v, w)$ is not a v-junction in $\tilde{G}$, where $v'$ is the predecessor of $v$. If $(v', v, u)$ is not a v-junction in $\tilde{G}$, then (D1) is true for $\gamma$ in $\tilde{G}$. If it is a v-junction, then $(v, u, w)$ is not and (D1) is true too.

Assume now that (D2) is true for $\gamma \setminus u$ in $\tilde{G}$. Again, there is no problem if (D2) occurs for some point other than $v$ or $w$, so let us consider the case for which it happens at $v$. This means that $v \notin \tilde{\mathcal{A}}_U$ and $(v', v, w)$ is a v-junction. But, since $u \in U$, the link between $u$ and $v$ must be from $u$ to $v$ in $\tilde{G}$ so that there is no v-junction at $u$ and (D1) is true in $\tilde{G}$. This proves lemma 15.20.

**Lemma 15.21** *Let $\gamma$ be a path with a three-point loop at $u \in U$ for G. Assume that $\gamma$ does not satisfy (D2) in G. Then $\gamma \setminus u$ does not satisfy this property either.*

Let us assume that $\gamma \setminus u$ satisfies (D2) and reach a contradiction. Letting $(v, u, w)$ be the three-point loop, (D2) can only happen in $\gamma \setminus u$ at $v$ or $w$, and let us assume that this happens at $v$, so that, $v'$ being the predecessor of $v$, $(v', v, w)$ is a v-junction in $G$ with $v \notin \mathcal{A}_U$. Since $v \notin \mathcal{A}_U$, the link between $u$ and $v$ in $G$ must be from $u$ to $v$, but this implies that $(v', v, u)$ is a v-junction in $G$ with $v \notin \mathcal{A}_U$ which is a contradiction: this proves lemma 15.21.

The previous three lemmas directly imply the next one.

**Lemma 15.22** *If $\gamma$ is a path between s and t that does not satisfy (D2) for G, then $\gamma$ satisfies (D1) or (D2) for $\tilde{G}$.*

Indeed, if we start with $\gamma$ that does not satisfy (D2) for $G$, lemma 15.21 allows us to progressively remove three-point loops from $\gamma$ until none remains with a final path that satisfies the assumptions of lemma 15.19 and therefore satisfies (D1) in $\tilde{G}$, and lemma 15.20 allows us to add the points that we have removed in reverse order while always satisfying (D1) or (D2) in $\tilde{G}$.

We now partially relax the hypothesis that (D2) is not satisfied with the next lemma.

**Lemma 15.23** *If $\gamma$ is a path between s and t that does not pass in $V \setminus \mathcal{A}_U$ at a linked v-junction for G, then $\gamma$ satisfies (D1) or (D2) for $\tilde{G}$.*

Assume that $\gamma$ does not satisfy (D2) for $\tilde{G}$ (otherwise the result is proved). By lemma 15.22, $\gamma$ must satisfy (D2) for $G$. So, take an intersection of $\gamma$ with $V \setminus \mathcal{A}_U$ that occurs at a v-junction in $G$, that we will denote $(v, u, w)$. This is still a v-junction in $\tilde{G}$ since we assume it to be unlinked. Since (D2) is false in $\tilde{G}$, we must have $u \in \tilde{\mathcal{A}}_U$, and there is an oriented path, $\tau$, from $u$ to $U$ in $\tilde{G}$.

We can assume that $\tau$ has no v-junction in $G$. If a v-junction exists in $\tau$, then this v-junction must be linked (otherwise this would also be a v-junction in $\tilde{G}$ and con-tradict the fact that $\tau$ is consistently oriented in $\tilde{G}$), and this link must be oriented from $u$ to $U$ in $\tilde{G}$ to avoid creating a loop in this graph. This implies that we can bypass the v-junction while keeping a consistently oriented path in $\tilde{G}$, and iterate this until $\tau$ has no v-junction in $G$. But this implies that $\tau$ is consistently oriented in $G$, necessarily from $U$ to $u$ since $u \notin \mathcal{A}_U$.

Denote $\tau = (u_0 = u, v_1, \ldots, u_n \in U)$. We now prove by induction that each $(v, u_k, w)$ is an unlinked v-junction. This is true when $k = 0$, and let us assume that it is true for $k - 1$. Then $(u_k, u_{k-1}, v)$ is a v-junction in $G$ but not in $\tilde{G}$: so it must be linked and there exists an edge between $v$ and $u_k$. In $\tilde{G}$, this edge must be oriented from $v$ to $u_k$, since $(v, u_{k-1}, u_k, v)$ would form a loop otherwise. For the same reason, there must be an edge in $\tilde{G}$ from $w$ to $u_k$ so that $(v, u_k, w)$ is an unlinked v-junction.

Since this is true for $k = n$, we can replace $u$ by $u_n$ in $\gamma$ and still obtain a valid path. This can be done for all intersections of $\gamma$ with $V \setminus \mathcal{A}_U$ that occur at v-junctions. This finally yields a path (denote it $\bar{\gamma}$) which does not satisfy (D2) in $G$ anymore, and therefore satisfies (D1) or (D2) in $\tilde{G}$: so $\bar{\gamma}$ must either pass in $U$ without a v-junction or in $V \setminus \tilde{A}_U$ at a v-junction. None of the nodes that were modified can satisfy any of these conditions, since they were all in $U$ with a v-junction, so that the result is true for the original $\gamma$ also. This proves lemma 15.23.

So the only unsolved case is when $\gamma$ is allowed to pass in $V \setminus \mathcal{A}_U$ at linked v-junctions. We define an algorithm that removes them as follows. Let $\gamma_0 = \gamma$ and let $\gamma_k$ be the path after step $k$ of the algorithm. One passes from $\gamma_k$ to $\gamma_{k+1}$ as follows.

- If $\gamma_k$ has no linked v-junctions in $V \setminus \mathcal{A}_U$ for $G$, stop.
- Otherwise, pick such a v-junction and let $(v, u, w)$ be the three nodes involved in it.

  (i) If $v \in U, v' \notin U$ and $(v', v, u)$ is a v-junction in $\tilde{G}$, remove $v$ from $\gamma_k$ to define $\gamma_{k+1}$.
  (ii) Otherwise, if $w \in U, w' \notin U$ and $(u, w, w')$ is a v-junction in $\tilde{G}$, remove $w$ from $\gamma_k$ to define $\gamma_{k+1}$.
  (iii) Otherwise, remove $u$ from $\gamma_k$ to define $\gamma_{k+1}$.

None of the considered cases can disconnect the path. This is clear for case (iii) since $v$ and $w$ are linked. For case (i), note that, in $G$, $(v', v, u)$ cannot be a v-junction since $(v, u, w)$ is one. This implies that the v-junction in $\tilde{G}$ must be linked and that $v'$ and $u$ are connected.

The algorithm will stop at some point with some $\gamma_n$ that does not have any linked v-junction in $V \setminus \mathcal{A}_U$ anymore, which implies that (D1) or (D2) is true in $\tilde{G}$ for $\gamma_n$. To prove that this statement holds for $\gamma$, it suffices to show that if (D1) or (D2) is true in $\tilde{G}$ with $\gamma_{k+1}$, it must have been true with $\gamma_k$ at each step of the algorithm. So let's assume that $\gamma_{k+1}$ satisfies (D1) or (D2) in $\tilde{G}$.

First assume that we passed from $\gamma_k$ to $\gamma_{k+1}$ via case (iii). Assume that (D2) is true for $\gamma_{k+1}$, with as usual the only interesting case being when this occurs at $v$ or $w$. Assume it occurs at $v$ so that $(v', v, w)$ is a v-junction and $v \notin \tilde{\mathcal{A}}_U$. If $(v', v, u)$ is a v-junction, then (D2) is true with $\gamma_k$. Otherwise, there is an edge from $v$ to $u$ in $\tilde{G}$ which also implies an edge from $w$ to $u$ since $(v, u, w, v)$ would be a loop otherwise. So $(v, u, w)$ is a v-junction in $\tilde{G}$, and $u$ cannot be in $\tilde{\mathcal{A}}_U$ since its parent, $v$ would be in that set also. So (D2) is true in $\tilde{G}$. Now, assume that (D1) is true at $v$, so that $(v', v, w)$ is not a v-junction and $v \in U$. If $(v', v, u)$ is not a v-junction either, we are done, so assume the contrary. If $v' \in U$, then we cannot have a v-junction at $v'$ and (D1) is true. But $v' \notin U$ is not possible since this leads to case (i).

Now assume that we passed from $\gamma_k$ to $\gamma_{k+1}$ via case (i). Assume that (D1) is true for $\gamma_k$: this cannot be at $v'$ since $v' \notin U$, neither at $u$ since $u \notin \mathcal{A}_U$, so it will also be true for $\gamma_{k+1}$. The same statement holds with (D2) since $(v', v, u)$ is a v-junction in $\tilde{G}$ with $v \in U$ which implies that both $v'$ and $u$ are in $\tilde{\mathcal{A}}_U$. Case (ii) is obviously addressed similarly.

    With this, the proof of theorem 15.18 is complete.                        ∎

### 15.2.5  Probabilistic inference: Sum-prod algorithm

We now discuss the issue of using the sum-prod algorithm to compute marginal probabilities, $\mathbb{P}(X^{(s)} = x^{(s)})$ for $s \in V$ when $X$ is a Bayesian network on $G = (V, E)$. By definition, $\mathbb{P}(X = x)$ can be written in the form

$$\mathbb{P}(X = x) = \prod_{C \in \mathcal{C}} \varphi_C(x^{(C)})$$

where $\mathcal{C}$ contains all subsets $C_s := \{s\} \cup pa(s)$, $s \in V$. Marginal probabilities can therefore be computed easily when the factor graph associated to $\mathcal{C}$ is acyclic, according to proposition 14.13. However, because of the specific form of the $\varphi_C$'s (they are conditional probabilities), the sum-prod algorithm can be analyzed in more detail, and provide correct results even when the factor graph is not acyclic.

    The general rules for the sum-prod algorithm are

$$\begin{cases} m_{sC}(x^{(s)}) \leftarrow \displaystyle\prod_{\tilde{C}, s \in \tilde{C}, \tilde{C} \neq C} m_{\tilde{C}s}(x^{(s)}) \\[2ex] m_{Cs}(x^{(s)}) \leftarrow \displaystyle\sum_{y^{(C)}: y^{(s)} = x^{(s)}} \varphi_C(y^{(C)}) \prod_{t \in C \setminus \{s\}} m_{tC}(y^{(t)}) \end{cases}$$

They take a particular form for Bayesian networks, using the fact that a vertex $s$ belongs to $C_s$, and to all $C_t$ for $t \in ch(s)$.

$$m_{sC_s}(x^{(s)}) \quad \leftarrow \quad \prod_{t \in ch(s)} m_{C_t s}(x^{(s)}),$$

$$m_{sC_t}(x^{(s)}) \quad \leftarrow \quad m_{C_s s}(x^{(s)}) \prod_{u \in ch(s), u \neq t} m_{C_u s}(x^{(s)}), \text{ for } t \in ch(s),$$

$$m_{C_s s}(x^{(s)}) \quad \leftarrow \quad \sum_{y^{(C_s)}, y^{(s)} = x^{(s)}} p_s(y^{(pa(s))}, x^{(s)}) \prod_{t \in pa(s)} m_{t C_s}(y^{(t)}),$$

$$m_{C_t s}(x^{(s)}) \quad \leftarrow \quad \sum_{y^{(C_t)}, y^{(s)} = x^{(s)}} p_t(x^{(s)} \wedge y^{(pa(s) \backslash \{t\})}, y^{(t)}) m_{t C_t}(y^{(t)}) \prod_{u \in pa(s), u \neq t} m_{u C_t}(y^{(u)}),$$

$$\text{for } t \in ch(s).$$

These relations imply that, if $pa(s) = \emptyset$ ($s$ is a root), then $m_{C_s s} = p_s(x^{(s)})$. Also, if $ch(s) = \emptyset$ ($s$ is a leaf) then $m_{sC_s} = 1$. The following proposition shows that many of the messages become constant over time.

**Proposition 15.24** *All upward messages, $m_{sC_s}$ and $m_{C_t s}$ with $t \in ch(s)$ become constant (independent from $x^{(s)}$) in finite time.*

PROOF This can be shown recursively as follows. Assume that, for a given $s$, $m_{t C_t}$ is constant for all $t \in ch(s)$ (this is true if $s$ is a leaf). Then,

$$
\begin{aligned}
m_{C_t s}(x^{(s)}) \quad \leftarrow \quad & \sum_{pey C_t, y^{(s)} = x^{(s)}} p_t(x^{(s)} \wedge y^{(pa(s) \backslash \{t\})}, y^{(t)}) m_{t C_t}(y^{(t)}) \prod_{u \in pa(s), u \neq t} m_{u C_t}(y^{(u)}), \\
= \quad & m_{t C_t} \sum_{y^{(C_t)}, y^{(s)} = x^{(s)}} p_t(x^{(s)} \wedge y^{(pa(s) \backslash \{t\})}, y^{(t)}) \prod_{u \in pa(s), u \neq t} m_{u C_t}(y^{(u)}) \\
= \quad & m_{t C_t} \sum_{y^{(C_t \backslash \{t\})}, y^{(s)} = x^{(s)}} \prod_{u \in pa(s), u \neq t} m_{u C_t}(y^{(u)}) \\
= \quad & m_{t C_t} \prod_{u \in pa(s), u \neq t} \sum_{y^{(u)}} m_{u C_t}(y^{(u)})
\end{aligned}
$$

which is constant. Now

$$m_{sC_s}(x^{(s)}) \leftarrow \prod_{t \in ch(s)} m_{C_t s}(x^{(s)})$$

is also constant. This proves that all $m_{sC_s}$ progressively become constant, and, as we have just seen, this implies the same property for $m_{C_t s}$, $t \in ch(s)$. ∎

This proposition implies that, if initialized with constant messages (or after a finite time), the sum-prod algorithm iterates

$$m_{sC_s} \leftarrow \prod_{t \in ch(s)} m_{C_t s}$$

$$m_{C_s s}(x^{(s)}) \leftarrow \sum_{y^{(C_s)}, y^{(s)} = x^{(s)}} p_s(y^{()} pa(s), x^{(s)}) \prod_{t \in pa(s)} m_{t C_s}(y^{(t)})$$

$$m_{sC_t}(x^{(s)}) \leftarrow m_{C_s s}(x^{(s)}) \prod_{u \in ch(s), u \neq t} m_{C_u s}, \quad t \in ch(s)$$

$$m_{C_t s} \leftarrow m_{sC_s} \prod_{u \in pa(t), u \neq s} \sum_{y^{(u)}} m_{u C_s}(y^{(u)}), \quad t \in ch(s).$$

From this expression, we can conclude

**Proposition 15.25** *If the previous algorithm is first initialized with upward messages, $m_{sC_s} = m_{C_t s}$ all equal to 1, and if downward messages are computed top down from the roots to the leaves, the obtained configuration of messages is invariant for the sum-prod algorithm.*

PROOF If all upward messages are equal to 1, then clearly, the downward messages sum to 1 once they are updated from roots to leaves, and this implies that the upward messages will remain equal to 1 for the next round. The obtained configuration is invariant since the downward messages are recursively uniquely defined by their value at the roots. ∎

The downward messages, under the previous assumptions, satisfy $m_{sC_t}(x^{(s)}) = m_{C_s s}(x^{(s)})$ for all $t \in ch(s)$ and therefore

$$m_{C_s s}(x^{(s)}) = \sum_{y^{(C_s)}, y^{(s)} = x^{(s)}} \pi(y^{(pa(s))}, x^{(s)}) \prod_{t \in pa(s)} m_{C_t t}(y^{(t)}). \tag{15.5}$$

Note that the associated "marginals" inferred by the sum-prod algorithm are

$$\sigma_s(x^{(s)}) = \prod_{C, s \in C} m_{Cs}(x^{(s)}) = m_{C_s s}(x^{(s)})$$

since $m_{C_t s}(x^{(s)}) = 1$ when $t \in ch(s)$.

Although the sum-prod algorithm initialized with unit messages converges to a stable configuration if run top-down, the obtained $\sigma_s$'s do not necessarily provide the correct single site marginals. There is a situation for which this is true, however, which is when the initial directed graph is singly connected, as we will see below.

Before this, let us analyze the complexity resulting from an iterative computation of the marginal probabilities, similar to what we have done with trees.

We define the depth of a vertex in $G$ as follows.

**Definition 15.26** *Let $G = (V, E)$ be a DAG. The depth of a vertex $s$ in $V$ is defined recursively by*

- depth$(s) = 0$ *if $s$ has no parent.*
- depth$(s) = 1 + \max\big(\text{depth}(t), t \in pa(s)\big)$ *otherwise.*

The recursive computation of marginal distributions is made possible (although not always feasible) with the following remark.

**Lemma 15.27** *Let $X$ be a Bayesian network on the DAG $G = (V, E)$, and $S \subset V$, such that all elements in $S$ have the same depth. Let $pa(S)$ be the set of parents of elements in $S$, and $T = \text{depth}^-(S)$ the set of vertexes in $V$ with depth strictly smaller than the depth of $S$. Then $(X^{(S)} \perp\!\!\!\perp X^{(T \setminus pa(S))} \mid X^{(pa(S))})$ and the variables $X^{(s)}, s \in S$ are conditionally independent given $X^{(pa(S))}$.*

PROOF It suffices to show that vertexes in $S$ are separated from $T \setminus pa(S)$ and from other elements of $S$ by $pa(S)$ for the graph $(G_{S \cup T})^{\sharp}$. Any path starting at $s \in S$ must either pass by a parent of $s$ (which is what we want), or by one of its children, or by another vertex that shares a child with $s$ in $G_{S \cup T}$. But $s$ cannot have any child in $G_{S \cup T}$, since this child cannot have a smaller depth than $s$, and it cannot be in $S$ either since all elements in $S$ have the same depth. $\blacksquare$

This lemma allows us to work recursively as follows. Assume that we can compute marginal distributions over sets $S$ with maximal depth no larger than $d$. Take a set $S$ of maximal depth $d + 1$, and let $S_0$ be the set of elements of depth $d + 1$ in $S$. Then, letting $T = \text{depth}^-(S) = \text{depth}^-(S_0)$, and $S_1 = S \setminus S_0$,

$$
\begin{aligned}
\mathbb{P}(X^{(S)} = x^{(S)}) &= \sum_{y^{(T \setminus S_1)}} \mathbb{P}(X^{(S_0)} = x^{(S_0)} \mid X^{(T)} = y^{(T \setminus S_1)} \wedge x^{(S_1)}) P(X^{(T \cup S_1)} = y^{(T \setminus S_1)} \wedge x^{(S_1)}) \\
&= \sum_{y^{(pa(S) \setminus S_1)}} \prod_{s \in S_0} p_s((y \wedge x)^{(pa(s))} \wedge x^{(S_1)}, x^{(s)}) \mathbb{P}(X^{(pa(S_0) \cup S_1)} = y^{(pa(S_0) \setminus S_1)} \wedge x^{(S_0)}) \quad (15.6)
\end{aligned}
$$

Since $pa(S) \cup S_1$ has maximal depth strictly smaller than the maximal depth of $S$, this indeed provides a recursive formula for the computation of marginal over subsets of $V$ with increasing maximal depths. However, because one needs to add parents to the considered set when reducing the depth, one may end up having to compute marginals over very large sets, which becomes intractable without further assumptions.

A way to reduce the complexity is to assume that the graph $G$ is singly connected, as defined below.

**Definition 15.28** *A DAG G is singly connected if there exists at most one path in G that connects any two vertexes.*

Such a property is true for a tree, but also holds for some networks with multiple parents. We have the following nice property in this case.

**Proposition 15.29** *Let G be a singly connected DAG and X a Bayesian network on G. If s is a vertex in G, the variables $(X^{(t)}, t \in pa(s))$ are mutually independent.*

PROOF We have, using proposition 15.5,

$$\mathbb{P}(X^{(pa(s))} = x^{(pa(s))}) = \sum_{y^{(\mathcal{A}_{pa(s)})}, y^{(pa(s))} = x^{(pa(s))}} \prod_{u \in \mathcal{A}_{pa(s)}} p_u(y^{(pa(u))}, y^{(u)}).$$

Because the graph is singly connected, two parents of $s$ cannot have a common ancestor (since there would then be two paths from this ancestor to $S$). So $\mathcal{A}_{pa(s)}$ is the disjoint union of the $\mathcal{A}_t$'s for $t \in pa(s)$ and we can write

$$
\begin{aligned}
\mathbb{P}(X^{(pa(s))} = x^{(pa(s))}) &= \sum_{y^{(\mathcal{A}_{pa(s)})}, y^{(pa(s))} = x^{(pa(s))}} \prod_{t \in pa(s)} \prod_{u \in \mathcal{A}_t} p_u(y^{(pa(u))}, y^{(u)}) \\
&= \prod_{t \in pa(s)} \sum_{y^{(\mathcal{A}_t)}, y^{(t)} = x^{(t)}} \prod_{u \in \mathcal{A}_t} p_u(y^{(pa(u))}, y^{(u)}) \\
&= \prod_{t \in pa(s)} \mathbb{P}(X^{(t)} = x^{(t)})
\end{aligned}
$$

This proves the lemma.                                                                                     ∎

Section 15.2.5 can be simplified under the assumption of a singly connected graph, at least for the computation of single vertex marginals; we have, if $s \in V$ and $G$ is singly connected

$$\mathbb{P}(X^{(s)} = x^{(s)}) = \sum_{y^{(pa(s))}} p_s(y^{(pa(s))}, x^{(s)}) \prod_{t \in pa(s)} \mathbb{P}(X^{(t)} = y^{(t)}). \tag{15.7}$$

This is now recursive in single vertex marginal probabilities. It moreover coincides with the recursive equation that defines the messages $m_{C_s s}$ in (15.5), which shows that the sum-prod algorithm provides the correct answer in this case.

### 15.2.6 Conditional probabilities and interventions

One of the main interests of graphical models is to provide an ability to infer the behavior of hidden variables of interest given other, observed, variables. When dealing with oriented graphs the way this should be analyzed is, however, ambiguous.

Let's consider an example, provided by the graph in fig. 15.1. The Bayesian net-



Figure 15.1: Example of causal graph.

work interpretation of this graph is that both events (which may be true or false) "Bad weather" and "Broken HVAC" happen first, and that they are independent. Then, given their observation, the "No school" event may occur, probably more likely if the weather is bad or the HVAC is broken or snow, and even more likely if both happened at the same time.

Now consider the following passive observation: you wake up, you haven't checked the weather yet or the news yet, and someone tells you that there is no school today. Then you may infer that there is more chances than usual for bad weather or the HVAC broken at school. Conditionally to this information, these two events become correlated, even if they were initially independent. So, even if the "No school" event is considered as a probabilistic consequence of its parents, observing it influences our knowledge on them.

Now, here is an intervention, or manipulation: the school superintendent has declared that he has given enough snow days for the year and declared that there would be school today whatever happens. So you know that the "no-school" event will not happen. Does it change the risk of bad weather of broken HVAC? Obviously not: an intervention on a node does not affect the distribution of the parents.

Manipulation and passive observation are two very different ways of affecting unobserved variables in Bayesian networks. Both of them may be relevant in applications. Of the two, the simplest to analyze is intervention, since it merely consists in clamping one of the variables while letting the rest of the network dynamics unchanged. This leads to the following formal definition of manipulation.

**Definition 15.30** *Let $G = (V, E)$ be a directed acyclic graph and $X$ a Bayesian network on*

*G. Let S be a subset of G and $x^{(S)} \in F_S$ a given configuration on S. Then the manipulated distribution of X with fixed values $x^{(S)}$ on S is the Bayesian network on the restricted graph $G_S$, with the same conditional probabilities, using the value $x^{(s)}$ every time a vertex $s \in S$ is a parent of $t \in V \setminus S$ in G.*

So, if the distribution of $X$ is given by (15.2), then its distribution after manipulation on $S$ is

$$\tilde{\pi}(y^{(V \setminus S)}) = \prod_{t \in V \setminus S} p_t(y^{(pa(t))}, y^{(t)})$$

where $pa(t)$ is the set of parents of $t$ in $G$, and $y^{(s)} = x^{(s)}$ whenever $s \in pa(t) \cap S$.

The distribution of a Bayesian network $X$ after passive observation $X^{(S)} = x^{(S)}$ is not so easily described. It is obviously the conditional distribution $P(X^{(V \setminus S)} = y^{(V \setminus S)} \mid X^{(S)} = x^{(S)})$ and therefore requires using the conditional dependency structure, involving the moral graph and/or d-separation.

Let us discuss this first in the simpler case of trees, for which the moral graph is the undirected acyclic graph underlying the tree, and d-separation is simple separation on this acyclic graph. We can then use proposition 13.22 to understand the new structure after conditioning: it is a $G^\flat_{V \setminus S}$-Markov random field, and, for $t \in V \setminus S$, the conditional distribution of $X^{(t)} = y^{(t)}$ given its neighbors is the same as before, using the value $x^{(s)}$ when $s \in S$. But note that when doing this (passing to $G^\flat$), we broke the causality relation between the variables. We can however always go back to a tree (or forest, since connectedness may have been broken) with the same edge orientation as they initially were, but this requires reconstituting the edge joint probabilities from the new acyclic graph, and therefore using (acyclic) belief propagation.

With general Bayesian networks, we know that the moral graph can be loopy and therefore a source of difficulties. The following proposition states that the damage is circumscribed to the ancestors of $S$.

**Proposition 15.31** *Let $G = (V, E)$ be a directed acyclic graph, X a Bayesian network on G, $S \subset V$ and $x^{(\mathcal{A}_S)} \in \mathcal{F}(\mathcal{A}_S)$. Then the conditional distribution of $X^{(\mathcal{A}_S^c)}$ given by $X^{(\mathcal{A}_S)} = x^{(\mathcal{A}_S)}$ coincides with the manipulated distribution in definition 15.30.*

PROOF  The conditional distribution is proportional to

$$\prod_{s \in V} p(y^{(pa(s))}, y^{(s)})$$

with $y^{(t)} = x^{(t)}$ if $t \in \mathcal{A}_S$. Since $s \in \mathcal{A}_S$ implies $pa(s) \subset \mathcal{A}_S$, all terms with $s \in \mathcal{A}_S$ are constant in the sum and can be factored out after normalization. So the conditional

distribution is proportional to

$$\prod_{s \in \mathcal{A}_S^c} p(y^{(pa(s))}, y^{(s)})$$

with $y^{(t)} = x^{(t)}$ if $t \in \mathcal{A}_S$. But we know that such products sum to 1, so that the conditional distribution is equal to this expression and therefore provides a Bayesian network on $G_{\mathcal{A}_S^c}$. ∎

## 15.3 Structural equation models

Structural equation models (SEM's) provides an alternative (and essentially equivalent) formulation of Bayesian networks, which may be more convenient to use, especially when dealing with variables taking values in general state spaces.

Let $G = (V, E)$ be a directed acyclic graph. SEMs are associated to families of functions $\Phi_s : \mathcal{F}(pa(s)) \times \mathcal{B}_s \to F_s$ and random variables $\xi_s : \Omega \to \mathcal{B}_s$ (where $\mathcal{B}_s$ is some measurable set), for $s \in V$. The random field $X : \Omega \to \mathcal{F}(V)$ associated to the SEM satisfies the equations

$$X_s = \Phi^{(s)}(X^{(s-)}, \xi^{(s)}). \tag{15.8}$$

Because of the DAG structure, these equations uniquely define $X$ once $\xi$ is specified. As a consequence, there exists a function $\Psi$ such that $X = \Psi(\xi)$.

The model is therefore fully specified by the functions $\Phi^{(s)}$ and the probability distributions of the variables $\xi^{(s)}$. We will assume that they have a density, denoted $g^{(s)}, s \in V$, with respect to some measure $\mu_s$ on $\mathcal{B}_s$. They are typically chosen as uniform distributions on $\mathcal{B}_s$ (continuous and compact, or discrete) or as standard Gaussian when $\mathcal{B}_s = \mathbb{R}^{d_s}$ for some $d_s$. One also generally assumes that the variables $(\xi^{(s)}, s \in V)$ are jointly independent, and we make this assumption below.

Let $V_k$, $k \geq 0$, be the set of vertexes in $V$ with depth $k$ (c.f. definition 15.26) and $V_{<k} = V_0 \cup \cdots \cup V_{k-1}$. Then (using the independence of $(\xi^{(s)}, s \in V)$), for $s \in V_k$, the conditional distribution of $X^{(s)}$ given $X^{(V_{<k})} = x^{(V_{<k})}$ is the distribution of $\Phi^{(s)}(x^{(s-)}, \xi^{(s)})$. Formally this is given by

$$\Phi^{(s)}(x^{(s-)}, \cdot)_\sharp (g^{(s)} \mu_s),$$

the pushforward of the distribution of $\xi^{(s)}$ by $\Phi^{(s)}(x^{(s-)}, \cdot)$.

More concretely, assume that $\xi_s$ follows a uniform distribution on $\mathcal{B}_s = [0,1]^h$ for some $h$, and assume that $F_s$ is finite for all $s$. Then,

$$P(X^{(s)} = x^{(s)} \mid X^{(V_{<k})} = x^{(V_{<k})}) = \text{Volume}(U_s(x^{(pa(s))}, x^{(s)})) \stackrel{\Delta}{=} p_s(x^{(pa(s))}, x^{(s)})$$

where

$$U_s(x^{(pa(s))}, x^{(s)}) = \left\{ \xi \in [0,1]^h : \Phi^{(s)}(x^{(s-)}, \xi) = x^{(s)} \right\}.$$

Since variables $X^{(s)}$, $s \in V_k$ are conditionally independent given $X^{(V<k)}$, we find that $X$ decomposes as a Bayesian network over $G$,

$$P(X = x) = \prod_{s \in V} p_s(x^{(pa(s))}, x^{(s)}).$$

Similarly, if $F_s = \mathcal{B}_s = \mathbb{R}^{d_s}$, $\xi^{(s)} \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^{d_s}})$, and $\xi^{(s)} \mapsto \Phi_\theta^{(s)}(x^{(pa(s))}, \xi^{(s)})$ is invertible, with $C^1$ inverse $x^{(s)} \mapsto \Psi_\theta^{(s)}(x^{(pa(s))}, x^{(s)})$, then $X$ is a Bayesian network, with continuous variables, and, using the change of variable formula, the conditional distribution of $X^{(s)}$ given $X^{(pa(s))} = x^{(s-)}$ has p.d.f.

$$p_s(x^{(pa(s))}, x^{(s)}) = \frac{1}{(2\pi)^{d_s/2}} \exp\left(-\frac{1}{2}|x_s - \Psi_\theta^{(s)}(x^{(pa(s))}, x^{(s)})|^2\right) \left|\det(\partial_{x^{(s)}} \Psi_\theta^{(s)}(x^{(pa(s))}, x^{(s)}))\right|.$$

A simple and commonly used special case for this example are linear SEMs, with

$$X^{(s)} = a_s + b_s^T X^{(s)} + \sigma_s \xi^{(s)}.$$

In this case, the inverse mapping is immediate and the Jacobian determinant in the change of variables is $1/\sigma_s^{d_s}$.

# Chapter 16

# Latent Variables and Variational Methods

## 16.1 Introduction

We will describe, in the next chapters, methods that fit a parametric model to the observation while introducing unobserved, or "latent," components in their models, whose inference typically attaches interpretable information or structure to the data. We have seen one such example in the form of the mixture of Gaussian in chapter 4, that we will revisit in chapter 19. We now provide a presentation of the variational Bayes paradigm that provides a general strategy to address latent variable problems [143, 97, 14, 100].

The general framework is as follows. Variables in the model are divided in two groups: the observable part, that we denote $X$, and the latent part, denoted $Z$. In many models $Z$ represents some unobservable structure, such that $X$ conditional to $Z$ has some relatively simple distribution (in a Bayesian estimation context, $Z$ often contains model parameters). The quantity of interest, however, is the conditional distribution of $Z$ given $X$ (also called the "posterior distribution"), which allows one to infer the latent structure from the observations, and will also have an important role in maximum likelihood parametric estimation, as we will see below. This conditional distribution is not always easy to compute or simulate, and variational Bayes provides a framework under which it can be approximated.

## 16.2 Variational principle

We consider a pair of random variables $X$ and $Z$, where $X$ is considered as "observed" and $Z$ is hidden, or "latent". We will use $U = (X, Z)$ to denote the two variables taken together. We denote as usual by $P_U$ the probability law of $U$, defined on $\mathcal{R}_U = \mathcal{R}_X \times \mathcal{R}_Z$ by $P_U(A) = \mathbb{P}(U \in A)$. We will also assume that there exists a measure $\mu$ on $\mathcal{R}_U$ that decomposes as a product measure $\mu = \mu_X \times \mu_Z$ (where $\mu_X$ and $\mu_Z$

are measures on $\mathcal{R}_X$ and $\mathcal{R}_Z$), such that $P_U \ll \mu$ ($\pi_U$ is absolutely continuous with respect to $\mu$). This implies that $P_U$ has a density with respect to $\mu$ that we will denote $f_U$. If both $\mathcal{R}_X$ and $\mathcal{R}_Z$ are discrete, $\mu$ is typically the counting measure, and if they are both Euclidean space, $\mu$ can be the Lebesgue measure on the product.[1]

The variables $X$ and $Z$ then have probability density functions with respect to $\mu_X$ ad $\mu_Z$, given by

$$f_X(x) = \int_{\mathcal{R}_Z} f_U(x,z)\mu_Z(dz) \quad \text{and} \quad f_Z(z) = \int_{\mathcal{R}_X} f_U(x,z)\mu_X(dx).$$

The conditional distribution of $X$ given $Z = z$, denoted $P_X(\cdot \mid Z = z)$, has density $f_X(x \mid z) = f_U(x,z)/f_Z(z)$ with respect to $\mu_X$ and that of $Z$ given $X = x$, denoted $P_Z(\cdot \mid X = x)$, has density $f_Z(z \mid x) = f_U(x,z)/f_X(x)$ with respect to $\mu_Z$. We will be mainly interested by approximations of $P_Z(\cdot \mid X = x)$, assuming that $P_Z$ and $P_X(\cdot \mid Z = z)$ (and hence $P_U$) are easy to compute or simulate.

We will use the Kullback-Liebler divergence to quantify the accuracy of the approximation. As stated in proposition 4.1, we have

$$P_Z(\cdot \mid X = x) = \underset{\nu \in \mathcal{M}_1(\mathcal{R}_Z)}{\operatorname{argmin}} \ KL(\nu \| P_Z(\cdot \mid X = x))$$

where $\mathcal{M}_1(\mathcal{R}_Z)$ denotes the set of all probability distributions on $\mathcal{R}_Z$. Note that all distributions $\nu$ for which $KL(\nu \| \pi_Z(\cdot | X = x))$ is finite must be absolutely continuous with respect to $\mu_Z$ and therefore take the form $\nu = g\mu_Z$. One has

$$KL(gd\mu_Z \| P_Z(\cdot | X = x)) = \int_{\mathcal{R}_Z} \log \frac{g(z)}{f_Z(z|x)} g(z)\mu_Z(dz)$$

$$= \int_{\mathcal{R}_Z} \log \frac{g(z)}{f_U(x,z)} g(z)\mu_Z(dz) + \log f_X(x). \quad (16.1)$$

We will denote by $\mathcal{P}(\mu_Z)$, or just $\mathcal{P}$ when there is no ambiguity, the set of all p.d.f.'s $g$ with respect to $\mu_Z$, i.e., the set of all non-negative measurable functions on $\mathcal{R}_Z$ with $\int_{\mathcal{R}_Z} g(z)\mu_Z(dz) = 1$.

The basic principle of variational Bayes methods is to replace $\mathcal{P}$ by a subset $\widehat{\mathcal{P}}$ and to define the approximation

$$\widehat{P_Z}(\cdot | X = x) = \underset{g \in \widehat{\mathcal{P}}}{\operatorname{argmin}} KL(g\mu_Z \| P_Z(\cdot | X = x)).$$

---

[1]The reader unfamiliar with measure theory may want to read this discussion by replacing $d\mu_X$ by $dx$, $d\mu_Z$ by dz and $d\mu_U$ by $dx\,dz$, i.e., in the context of continuous probability distributions having p.d.f.'s with respect to the Lebesgue's measure.

For the approximation to be practical, the set $\widehat{\mathcal{P}}$ must obviously be chosen so that the computation of $\widehat{P_Z}(\cdot|X = x)$ is computationally feasible. We now review a few examples, before passing to the EM algorithm and its approximations.

## 16.3 Examples

### 16.3.1 Mode approximation

Assume that $\mathcal{R}_Z$ is discrete and $\mu_Z$ is the counting measure so that

$$KL(g\,d\mu_Z\,\|\,P_Z(\cdot\mid X = x)) - \log f_X(x) = \sum_{z \in \mathcal{R}_Z} \log \frac{g(z)}{f_U(x,z)} g(z),$$

the sum being infinite if there exists $z$ such that $\nu(z) > 0$ and $f_U(x,z) = 0$. Take

$$\widehat{\mathcal{P}} = \{\mathbf{1}_z : z \in \mathcal{R}_Z\},$$

the family of all Dirac functions on $\mathcal{R}_Z$. Then,

$$KL(\mathbf{1}_z\,\|\,P_Z(\cdot|X = x)) - \log f_X(x) = -\log f_U(x,z).$$

The variational approximation of $P_Z(\cdot\mid X = x)$ over $\widehat{\mathcal{P}}$ therefore is the Dirac measure at point(s) $z \in \mathcal{R}_Z$ at which $f_U(x,z)$ is largest, i.e., the mode(s) of the posterior distribution. This approximation is often called the MAP approximation (for maximum a posteriori).

If $\mathcal{R}_Z$ is, say, $\mathbb{R}^q$ and $\mu_Z = dz$ is Lebesgue's measure, then the previous construction does not work because $\mathbf{1}_z$ is not a p.d.f. with respect to $\mu_Z$. In place of Dirac functions, one can use constant functions on small balls. Let $B(z,\epsilon)$ denote the open ball with radius $\epsilon$, and let $|B(z,\epsilon)|$ denote its volume. Let $\mathfrak{u}_{z,\epsilon} = \mathbf{1}_{B(z,\epsilon)}/|B(z,\epsilon)|$. Fixing $\epsilon$, we can consider the set

$$\widehat{\mathcal{P}} = \{\mathfrak{u}_{z,\epsilon} : z \in \mathbb{R}^q\}.$$

Now, one has (leaving the computation to the reader)

$$KL(\mathfrak{u}_{z,\epsilon}\,dz\,\|\,P_Z(\cdot|X = x)) - \log f_X(x) = -\log\left(\frac{1}{|B(z,\epsilon)|}\int_{B(z,\epsilon)} f_U(x,z')dz'\right).$$

The limit for small $\epsilon$ (assuming that $f_U(x,\cdot)$ is continuous at $z$, or defining the limit up to sets of measure zero) is $-\log f_U(x,z)$, justifying again choosing the mode of the posterior distribution of $Z$ for the approximation.

The mode approximation has some limitations. First, it is in general a very crude approximation of the posterior distribution. Second, even with the assumption that $f_U$ has closed form, this p.d.f. is often difficult to maximize (for example when defining models over large discrete sets). In such cases, the mode approximation has limited practical use.

### 16.3.2   Gaussian approximation

Let us still assume that $\mathcal{R}_Z = \mathbb{R}^q$ and that $\mu_Z = dz$. Let $\widehat{\mathcal{P}}$ be the family of all Gaussian distributions $\mathcal{N}(m, \Sigma)$ on $\mathbb{R}^q$. Then, denoting by $\varphi(\cdot; m, \Sigma)$ the density of $\mathcal{N}(m, \Sigma)$,

$$KL(\varphi(\cdot; m, \Sigma)\|P_Z(\cdot \mid X = x)) - \log f_X(x) = -\frac{q}{2}\log 2\pi - \frac{q}{2} - \frac{1}{2}\log \det(\Sigma)$$
$$- \int_{\mathbb{R}^q} \log f_U(x, z)\varphi(z; m, \Sigma)dz.$$

In order to provide the best approximation, $m$ and $\Sigma$ must therefore maximize

$$\int_{\mathbb{R}^q} \log f_U(x, z)\varphi(z; m, \Sigma)dz + \frac{1}{2}\log \det(\Sigma). \tag{16.2}$$

The resulting optimization problem does not have a closed form solution in general (see section 18.2.2 for an example in which stochastic gradient methods are used to solve this problem). Another approach that is commonly used in practice is to push the approximation further by replacing $\log f_U(x, z)$ by its second order expansion around its maximum as a function of $z$. Let $m(x)$ be the posterior mode, i.e., the value of $z$ at which $x \mapsto \log f_U(x, z)$ is maximal, that we will assume to be unique. Let $H(x)$ denote the $q \times q$ Hessian matrix formed by the second partial derivatives of $-\log f_U(x, z)$ (with respect to $z$) at $z = m(x)$. This matrix is positive semidefinite according to the choice made for $m(x)$, and we will assume that it is positive definite. Since the first derivatives of $\log f_U(x, z)$ at $m(x)$ must vanish, we have the expansion:

$$\log f_U(x, z) = \log f_U(x, m(x)) - \frac{1}{2}(z - m(x))^T H(x)(z - m(x)) + \cdots$$

Plugging the expansion into the integral in (16.2) yields

$$-\frac{1}{2}\text{trace}(H(x)\Sigma) - \frac{1}{2}(m - m(x))^T H(x)(m - m(x)) + \frac{1}{2}\log \det \Sigma.$$

To maximize this expression, one must clearly take $m = m(x)$. Moreover,

$$\partial_\Sigma \left(-\text{trace}(H(x)\Sigma) + \log \det \Sigma\right) = -H(x)^T + (\Sigma^T)^{-1} = -H(x) + \Sigma^{-1},$$

and we see that one must take $\Sigma = H(x)^{-1}$. This provides the Laplace approximation [62] of the posterior, $\mathcal{N}(m(x), H(x)^{-1})$, which is practical when the mode and corresponding second derivatives are feasible to compute.

### 16.3.3   Mean-field approximation

This section generalizes the approach discussed in proposition 14.6 for Markov random fields. Assume that $\mathcal{R}_Z$ can be decomposed into several components $\mathcal{R}_Z^{[1]}, \ldots, \mathcal{R}_Z^{[K]}$,

writing $z = (z^{[1]}, \ldots, z^{[K]})$ (for example, taking $K = q$ and $z^{[i]} = z^{(i)}$, the $i$th coordinate of $z$ if $\mathcal{R}_Z = \mathbb{R}^q$). Also assume that $\mu_Z$ splits into a product measure $\mu_Z^{[1]} \otimes \cdots \otimes \mu_Z^{[K]}$. Mean-field approximation consists in assuming that probabilities $\nu$ in $\widehat{\mathcal{P}}$ split into independent components, i.e., their densities $g$ take the form:

$$g(z) = g^{[1]}(z^{[1]}) \cdots g^{[K]}(z^{[K]}).$$

Then,

$$KL(\nu \| P_Z(\cdot \mid X = x)) - \log f_X(x) = \sum_{j=1}^{K} \int_{\mathcal{R}_Z^{[j]}} \log g^{[j]}(z^{[j]}) g^{[j]}(z^{[j]}) \mu_Z^{[j]}(dz^{[j]})$$

$$- \int_{\mathcal{R}_Z} \log f_U(x, z) \prod_{j=1}^{q} g^{[j]}(z^{[j]}) \mu_Z(dz). \quad (16.3)$$

The mean-field approximation may be feasible when $\log f_U(x, z)$ can be written as a sum of products of functions of each $z^{[j]}$. Indeed, assume that

$$\log f_U(x, z) = \sum_{\alpha \in A} \prod_{j=1}^{K} \psi_{\alpha, j}(z^{[j]}, x) \quad (16.4)$$

where $A$ is a finite set. To shorten notation, let us denote by $\langle \psi \rangle$ the expectation of a function $\psi$ with respect to the product p.d.f. $g$. Then, (16.3) can be written as

$$KL(\nu \| P_Z(\cdot \mid X = x)) - \log f_X(x) = \sum_{j=1}^{K} \langle \log g^{(j)}(z^{[j]}) \rangle - \sum_{\alpha \in A} \prod_{j=1}^{K} \langle \psi_{\alpha, j}(z^{[j]}, x) \rangle.$$

The following lemma will allow us to identify the form taken by the optimal p.d.f. $g^{[j]}$.

**Lemma 16.1** *Let $Q$ be a set equipped with a positive measure $\mu$. Let $\psi : Q \to \mathbb{R}$ be a measurable function such that*

$$C_\psi \overset{\Delta}{=} \int_Q \exp(\psi(q)) \mu(dq) < \infty.$$

*Let*

$$g_\psi(q) = \frac{1}{C_\psi} \exp(\psi(q)).$$

*Let $g$ be any p.d.f. with respect to $\mu$, and define*

$$F(g) = \int_Q (\log g(q) - \psi(q)) g(q) \mu(dq).$$

*Then $F(g_\psi) \leq F(g)$.*

PROOF  We note that $g_\psi > 0$, and that

$$KL(g\|g_\psi) = F(g) + \log C_\psi = F(g) - F(g_\psi),$$

which proves the result, since KL divergences are always non-negative.          ∎


Applying this lemma separately to each function $g^{[j]}$ implies that any optimal $g$ must be such that

$$g^{[j]}(z^{[j]}) \propto \exp\left(\sum_{\alpha \in A} M_{\alpha,j}\psi_{\alpha,j}(z^{[j]},x)\right)$$

with

$$M_{\alpha,j} = \prod_{j'=1,j'\neq j}^{K} \langle \psi_{\alpha,j'}(z^{[j']},x)\rangle.$$

We therefore have

$$\langle \psi_{\alpha,j}(z^{[j]},x)\rangle = \frac{\int_{\mathcal{R}_Z^{[j]}} \psi_{\alpha,j}(z^{[j]},x)\exp\left(\sum_{\alpha'\in A} M_{\alpha',j}\psi_{\alpha',j}(z^{[j]},x)\right)\mu_Z^{[j]}(dz^{[j]})}{\int_{\mathcal{R}_Z^{[j]}}\exp\left(\sum_{\alpha'\in A} M_{\alpha',j}\psi_{\alpha',j}(z^{[j]},x)\right)\mu_Z^{[j]}(dz^{[j]})} \tag{16.5}$$


This specifies a relationship expressing $\langle \psi_{\alpha,j}(z^{[j]},x)\rangle$ as a function of the other expectations $\langle \psi_{\alpha',j'}(z^{(j')},x)\rangle$ for $j \neq j'$. These equations put together are called the *mean-field consistency equations*. When these equations can be written explicitly, i.e., when the integrals in (16.5) can be evaluated analytically (which is generally the case when the p.d.f.'s $g^{[j]}$ can be associated with standard distributions), one obtains an algorithm that iterates (16.5) over all $\alpha$ and $j$ until stabilization (each step reducing the objective function in (16.3)).

Let us retrieve the result obtained in proposition 14.6 using the current formalism. Assume that $\mathcal{R}_X$ finite and $\mathcal{R}_Z = \{0,1\}^L$, where $L$ can be a large number, with

$$f_U(x,z) = \frac{1}{C}\exp\left(\sum_{j=1}^{L} \alpha_j(x)z^{(j)} + \sum_{i,j=1,i<j}^{L} \beta_{ij}(x)z^{(i)}z^{(j)}\right).$$

Take $K = L$, $z^{[j]} = z^{(j)}$. Applying the previous discussion, we see that $g^{[j]}$ must take the form

$$g^{[j]}(z^{(j)}) = \frac{\exp\left(\alpha_j(x)z^{(j)} + \sum_{i\neq j}\beta_{ij}(x)\langle z^{(i)}\rangle z^{(j)}\right)}{1 + \exp\left(\alpha_j(x) + \sum_{i\neq j}\beta_{ij}(x)\langle z^{(i)}\rangle\right)}$$

In particular

$$\langle z^{(j)}\rangle = \frac{\exp\left(\alpha_j(x) + \sum_{i\neq j}\beta_{ij}(x)\langle z^{(i)}\rangle\right)}{1 + \exp\left(\alpha_j(x) + \sum_{i\neq j}\beta_{ij}(x)\langle z^{(i)}\rangle\right)}$$

providing the mean-field consistency equations.

In this special case, it is also possible to express the objective function as a simple function of the expectations $\langle z^{(j)} \rangle$'s. We indeed have, letting $\rho_j = \langle z^{(j)} \rangle$,

$$\sum_{z \in \mathcal{R}_Z} \log f_U(x,z) \prod_{j=1}^{L} g^{[j]}(z^{(j)}) = -\log C + \sum_{j=1}^{L} \alpha_j(x)\rho_j + \sum_{i,j=1,i<j}^{L} \beta_{ij}(x)\rho_i \rho_j.$$

The values of $\rho_1, \ldots, \rho_L$ are then obtained by maximizing

$$\sum_{j=1}^{L} \alpha_j(x)\rho_j + \sum_{i,j=1,i<j}^{L} \beta_{ij}(x)\rho_i \rho_j - \sum_{j=1}^{L} \Big(\rho_j \log \rho_j + (1-\rho_j)\log(1-\rho_j)\Big).$$

The consistency equations express the fact that the derivatives of this expression with respect to each $\rho_j$ vanish.

## 16.4 Maximum likelihood estimation

### 16.4.1 The EM algorithm

We now consider maximum likelihood estimation with latent variables and use the notation of section 16.2. The main tool is the following obvious consequence of (16.1).

**Proposition 16.2** *One has*

$$\log f_X(x) = \max_{g \in \mathcal{P}(\mu_Z)} \int_{\mathcal{R}_Z} \log\left(\frac{f_U(x,z)}{g(z)}\right) g(z) d\mu_Z(z)$$

*and the maximum is achieved for $g(z) = f_Z(z \mid x)$, the conditional p.d.f. of $Z$ given $X = x$.*

Proof Equation (16.1) implies that

$$\int_{\mathcal{R}_Z} \log\left(\frac{f_U(x,z)}{g(z)}\right) g(z) d\mu_Z(z) = \log f_X(x) - KL(g\,\mu_Z \| P_Z(\cdot | X = x))$$

and the r.h.s. is indeed maximum when the Kullback-Liebler divergence vanishes, that is, when $g$ is the p.d.f. of $P_Z(\cdot \mid X = x)$. ∎

We will use this proposition for the derivation of the expectation-maximization (or EM) algorithm for maximum likelihood with latent variables. We now assume that $P_U$, and therefore $f_U$, is parametrized by $\theta \in \Theta$, and that a training set $T =$

$(x_1, \ldots, x_N)$ of $X$ is observed. To indicate the dependence in $\theta$, we will write $f_U(x, z; \theta)$, or $f_Z(z \mid x; \theta)$. The maximum likelihood estimator (m.l.e.) then maximizes

$$\ell(\theta) = \sum_{x \in T} \log f_X(x; \theta).$$

The EM algorithm is useful when the computation of the m.l.e. for complete observations, i.e., the maximization of

$$\log f_U(x, z; \theta)$$

when both $x$ and $z$ are given, is easy, whereas the same problem with the marginal distribution is hard.

From the proposition, we have:

$$\sum_{x \in T} \log f_X(x; \theta) = \sum_{x \in T} \max_{g_x \in \mathcal{P}(\mu_Z)} \int_{\mathcal{R}_Z} \log\left(\frac{f_U(x, z; \theta)}{g_x(z)}\right) g_x(z) \mu_Z(dz)$$

Therefore the maximum likelihood requires to compute

$$\max_{\theta, g_x, x \in T} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(\frac{f_U(x, z; \theta)}{g_x(z)}\right) g_x(z) \mu_Z(dz). \tag{16.6}$$

The maximization can therefore be done by iterating the following two steps.

1. Given $\theta_n$, compute

$$\operatorname*{argmax}_{g_x, x \in T} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(\frac{f_U(x, z; \theta)}{g_x(z)}\right) g_x(z) \mu_Z(dz).$$

2. Given $g_1, \ldots, g_N$, compute

$$\operatorname*{argmax}_{\theta} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(\frac{f_U(x, z; \theta)}{g_x(z)}\right) g_x(z) \mu_Z(dz)$$

$$= \operatorname*{argmax}_{\theta} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta)\right) g_x(z) \mu_Z(dz).$$

Step 1. is explicit and its solution is $g_x(z) = f_Z(z \mid x; \theta)$. Using this, both steps can be grouped together, yielding the EM algorithm.

**Algorithm 16.1 (EM algorithm)**
Let a statistical model with density $f_U(x, z; \theta)$ modeling an observable variable $X$ and a latent variable $Z$ be given, and a training set $T = (x_1, \ldots, x_N)$. Starting with an initial guess of the parameter, $\theta(0)$, the EM algorithm iterate the following equation until numerical stabilization, .

$$\theta_{n+1} = \operatorname*{argmax}_{\theta'} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta')\right) f_Z(z \mid x; \theta_n) \mu_Z(dz). \tag{16.7}$$

Equation (16.7) maximizes (in $\theta'$) a function defined as an expectation (for $\theta_n$), justifying the name "Expectation-Maximization."

### 16.4.2 Application: Mixtures of Gaussian

A mixture of Gaussian (MoG) model was introduced in chapter 4 ((4.4)). We now reinterpret it (in a slightly generalized version) as a model with partial observations and show how the EM algorithm can be applied. Let $\varphi(x; m, \Sigma)$ denote the p.d.f. of the $d$-dimensional multivariate Gaussian distribution with mean $m$ and covariance matrix $\Sigma$. We model $f_X(x; \theta)$ as

$$f_X(x; \theta) = \sum_{j=1}^{p} \alpha_j \varphi(x, ; c_j, \Sigma_j).$$

Here, $\theta$ contains all sequences $\alpha_1, \ldots, \alpha_p$ (non-negative numbers that sum to one), $c_1, \ldots, c_p \in \mathbb{R}^d$ and $\Sigma_1, \ldots, \Sigma_p$ ($d \times d$ positive definite matrices).

Using the previous notation, we therefore have $\mathcal{R}_X = \mathbb{R}^d$, and $\mu_X$ the Lebesgue measure on that space. The variable $Z$ will take values in $\mathcal{R}_Z = \{1, \ldots, p\}$, with $\mu_Z$ being the counting measure. We model the joint density function for $(X, Z)$ as

$$f_U(x, z; \theta) = \alpha_z \varphi(x; c_z, \Sigma_z). \tag{16.8}$$

Clearly $f_X$ is the marginal p.d.f. of $f_U$. One can therefore consider $Z$ as a latent variable, and therefore estimate $\theta$ using the EM algorithm.

We now make (16.7) explicit for mixtures of Gaussian. For given $\theta$ and $\theta'$ and $x \in \mathcal{R}$, let

$$U_x(\theta, \theta') = \frac{d}{2} \log 2\pi + \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta')\right) f_Z(z|x; \theta) d\mu_Z(z)$$

$$= \sum_{z=1}^{p} \left(\log \alpha'_z - \frac{1}{2} \log \det \Sigma'_z - \frac{1}{2}(x - c'_z)^T \Sigma'^{-1}_z (x - c'_z)\right) f_Z(z|x; \theta)$$

with

$$f_Z(z \mid x; \theta) = \frac{(\det \Sigma_z)^{-\frac{1}{2}} \alpha_z e^{-\frac{1}{2}(x - c_z)^T \Sigma_z^{-1}(x - c_z)}}{\sum_{j=1}^p (\det \Sigma_j)^{-\frac{1}{2}} \alpha_j e^{-\frac{1}{2}(x - c_j)^T \Sigma_j^{-1}(x - c_j)}}.$$

If $\theta_n$ is the current parameter in the EM, the next one, $\theta_{n+1}$ must maximize $\sum_{x \in T}^N U_x(\theta_n, \theta')$. This can be solved in closed form. To compute $\alpha'_1, \ldots, \alpha'_p$, one must maximize

$$\sum_{x \in T} \sum_{z=1}^p (\log \alpha'_z) f_Z(z|x; \theta)$$

subject to the constraint that $\sum_z \alpha'_z = 1$. This yields

$$\alpha'_z = \sum_{x \in T} f_Z(z|x; \theta) \Big/ \sum_{j=1}^p \sum_{x \in T} f_Z(j|x; \theta) = \zeta_z / N$$

with $\zeta_z = \sum_{x \in T} f_Z(z|x; \theta)$.

The centers $c'_1, \ldots, c'_p$ must minimize $\sum_{x \in T} (x - c'_z)^T \Sigma_z'^{-1}(x - c'_z) f_Z(z|x; \theta)$, which yields

$$c'_z = \frac{1}{\zeta_z} \sum_{x \in T} x f_Z(z|x; \theta).$$

Finally, $\Sigma'_z$ must minimize

$$\frac{\zeta_z}{2} \log \det \Sigma'_z + \frac{1}{2} \sum_{x \in T} (x - c'_z)^T \Sigma_z'^{-1}(x - c'_z) f_Z(z|x; \theta),$$

which yields

$$\Sigma'_z = \frac{1}{\zeta_z} \sum_{x \in T} (x - c'_z)(x - c'_z)^T f_Z(z|x; \theta).$$

We can now summarize the algorithm.

---

**Algorithm 16.2 (EM for Mixture of Gaussian distributions)**
   1.  Initialize the parameter $\theta(0) = (\alpha(0), c(0), \Sigma(0))$. Choose a small constant $\epsilon$ and a maximal number of iterations $M$.

   2.  At step $n$ of the algorithm, let $\theta = \theta(n)$ be the current parameter, writing for short $\theta = (\alpha, c, \Sigma)$.

   3.  Compute, for $x \in T$ and $i = 1, \ldots, p$

$$f_Z(i \mid x; \theta) = \frac{(\det \Sigma_i)^{-\frac{1}{2}} \alpha_i e^{-\frac{1}{2}(x - c_i)^T \Sigma_i^{-1}(x - c_i)}}{\sum_{j=1}^p (\det \Sigma_j)^{-\frac{1}{2}} \alpha_j e^{-\frac{1}{2}(x - c_j)^T \Sigma_j^{-1}(x - c_j)}}$$

and let $\zeta_i = \sum_{x \in T} f_Z(i|x; \theta)$, $i = 1, \ldots, p$.

4. Let $\alpha_i' = \zeta_i/N$.

5. For $i = 1, \ldots, p$, let

$$c_i' = \frac{1}{\zeta_i} \sum_{x \in T} x f_Z(i \mid x ; \theta).$$

6. For $i = 1, \ldots, p$, let

$$\Sigma_i' = \frac{1}{\zeta_i} \sum_{x \in T} (x - c_i')(x - c_i')^T f_Z(i \mid x ; \theta).$$

7. Let $\theta' = (\mu', c', \Sigma')$. If $|\theta' - \theta| < \epsilon$ or $n + 1 = M$: return $\theta'$ and exit the algorithm.

8. Set $\theta(n + 1) = \theta'$ and return to step 2.

---

**Remark 16.3** Algorithm 16.2 can be simplified by making restrictions on the model. Here are some examples.

(i) One may restrict to $\Sigma_i = \sigma_i^2 \mathrm{Id}_{\mathbb{R}^d}$ to reduce the number of free parameters. Then, Step 7 of the algorithm needs to be replaced by:

$$(\sigma_i')^2 = \frac{1}{d\zeta_i} \sum_{x \in T} |x - c_i'|^2 f_Z(i \mid x ; \theta).$$

(ii) Alternatively, the model may be simplified by assuming that all covariance matrices coincide: $\Sigma_i = \Sigma$ for $i = 1, \ldots, p$. Then, Step 7 becomes

$$\Sigma_i' = \frac{1}{N} \sum_{i=1}^{p} \sum_{x \in T} (x - c_i')(x - c_i')^T f_Z(i \mid x ; \theta). \qquad \blacklozenge$$

(iii) Finally, one may assume that $\Sigma$ is known and fixed in the algorithm (usually in the form $\Sigma = \sigma^2 \mathrm{Id}_{\mathbb{R}^d}$ for some $\sigma > 0$) so that Step 7 of the algorithm can be removed.

(iv) One may also assume also that the (prior) class probabilities are known, typically set to $\alpha_i = 1/p$ for all i, so that Step 4 can be skipped.

### 16.4.3 Stochastic approximation EM

The stochastic approximation EM (or SAEM) algorithm has been proposed by De-lyon et al. [58] (see this reference for convergence results) to address the situation in which the expectations for the posterior distribution cannot be computed in closed form, but can be estimated using Monte-Carlo simulations. SAEM uses a special

form of stochastic approximation, different from the SGD algorithm described in section 3.3. It updates, at each step $n$, an approximate objective function that we will denote $\lambda_n$ and a current parameter $\theta(n)$. It implements the following iterations:

$$
\begin{cases}
\xi_{n+1}^{(x)} \sim P_Z(\cdot \mid X = x; \theta_n), \quad x \in T \\
\lambda_{n+1}(\theta') = \left(1 - \frac{1}{n+1}\right)\lambda_n(\theta') + \frac{1}{n+1}\left(\sum_{x \in T} \log f_U(x, \xi_{n+1}^{(x)}; \theta') - \lambda_n(\theta')\right), \; \theta' \in \Theta \\
\theta_{n+1} = \underset{\theta'}{\operatorname{argmax}} \, \lambda_{n+1}(\theta')
\end{cases}
\tag{16.9}
$$

The second step means that

$$
\lambda_n(\theta') = \sum_{x \in T} \left( \frac{1}{n} \sum_{j=1}^{n} \log f_U(x, \xi_j^{(x)}; \theta') \right).
$$

Given that $\xi_{n+1}^{(x)} \sim P_Z(\cdot \mid X = x; \theta_n)$, one expects this expression to approximate

$$
\sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta')\right) f_Z(z \mid x; \theta) d\mu_Z(z)
$$

so that the third step of (16.9) can be seen as an approximation of (16.7). Sufficient conditions under which this actually happens (and $\theta(n)$ converges to a local maximizer of the likelihood) are provided in Delyon et al. [58] (see also Kuhn and Lavielle [112] for a convergence result under more general hypotheses on how $\xi$ is simulated).

   To be able to run this algorithm efficiently, one needs the simulation of the posterior distribution to be feasible. Importantly, one also needs to be able to update efficiently the function $\lambda_n$. This can be achieved when the considered model belongs to an exponential family, which corresponds to assuming that the p.d.f. of $U$ takes the form

$$
f_U(x, z; \theta) = \frac{1}{C(\theta)} \exp\left(\psi(\theta)^T H(x, z)\right)
$$

for some functions $\psi$ and $H$. For example, the MoG model of equation (4.4) takes

this form, with

$$\psi(\theta)^T = \Big(\quad \log\alpha_1 - \frac{1}{2}m_1^T\Sigma_1^{-1}m_1 - \frac{1}{2}\log\det\Sigma_1,\ldots,\log\alpha_p - \frac{1}{2}m_p^T\Sigma_p^{-1}m_p - \frac{1}{2}\log\det\Sigma_p,$$

$$\Sigma_1^{-1}m_1,\ldots,\Sigma_p^{-1}m_p,$$

$$\Sigma_1^{-1},\ldots,\Sigma_p^{-1}\Big),$$

$$H(x,z)^T = \Big(\quad \mathbf{1}_{z=1},\ldots,\mathbf{1}_{z=p},$$

$$x\mathbf{1}_{z=1},\ldots,x\mathbf{1}_{z=p},$$

$$-\frac{1}{2}xx^T\mathbf{1}_{z=1},\ldots,-\frac{1}{2}xx^T\mathbf{1}_{z=p}\Big)$$

and $C(\theta) = (2\pi)^{pd/2}$.

For such a model, we can replace the algorithm in (16.9) by the more manageable one:

$$\begin{cases} \xi_{n+1}^{(x)} \sim P_Z(\cdot \mid X = x; \theta_n), \ x \in T \\[2mm] \eta_{n+1}^{(x)} = \Big(1 - \frac{1}{n+1}\Big)\eta_n^{(x)} + \frac{1}{n+1}(H(x,\xi_{n+1}^{(x)}) - \eta_n^{(x)}) \\[2mm] \lambda_{n+1}(\theta') = \psi(\theta')^T\Big(\sum_{x\in T}\eta_{n+1}^{(x)}\Big) - \log C(\theta') \\[2mm] \theta_{n+1} = \underset{\theta'}{\mathrm{argmax}}\,\lambda_{n+1}(\theta') \end{cases} \tag{16.10}$$

We leave as an exercise the computation leading to the implementation of this algorithm for mixtures of Gaussian.

### 16.4.4 Variational approximation

Returning to proposition 16.2 and (16.6), we see that one can make a variational approximation of the maximum likelihood by computing

$$\max_{\theta\in\Theta, g_x\in\widehat{\mathcal{P}}, x\in T} \sum_{x\in T} \int_{\mathcal{R}_Z} \log\Big(\frac{f_U(x,z;\theta)}{g_x(z)}\Big) g_x(z)\mu_Z(dz), \tag{16.11}$$

where $\widehat{\mathcal{P}} \subset \mathcal{P}$ is a class of p.d.f. with respect to $\mu_Z$. The resulting algorithm is then implemented by iterating the computation of $g_x$, $x \in T$, using approximations similar to those provided in section 16.3, and maximization in $\theta$ for given $g_x$, $x \in T$. This variational approximation of the maximum likelihood estimator is therefore provided by the following algorithm.

**Algorithm 16.3 (Variational Bayes approximation of the m.l.e.)**
Let a statistical model with density $f_U(x,z;\theta)$ modeling an observable variable $X$ and a latent variable $Z$ be given, and a training set $T = (x_1,\dots,x_N)$ be observed. Let $\widehat{\mathcal{P}}$ be a set of p.d.f. on $\mathcal{R}_Z$ and define

$$\widehat{g}(\cdot\,;x,\theta) = \underset{g\in\widehat{\mathcal{P}}}{\mathrm{argmin}}\int_{\mathcal{R}_Z}\log\left(\frac{g(z)}{f_U(x,z;\theta)}\right)g(z)\mu_Z(dz)$$

(assuming that this minimizer is uniquely defined).

Starting with an initial guess of the parameter, $\theta_0$, iterate the following equation until numerical stabilization:

$$\theta(n+1) = \underset{\theta'}{\mathrm{argmax}}\sum_{x\in T}\int_{\mathcal{R}_Z}\log\left(f_U(x,z;\theta')\right)\widehat{g}(z|x;\theta(n))\mu_Z(dz). \tag{16.12}$$

---

Assume that the distributions in $\widehat{\mathcal{P}}$ are also parametrized, denoting their parameter by $\eta$, belonging to some Euclidean domain $H$. Let $g(\cdot;\eta)$ denote the p.d.f. in $\widehat{\mathcal{P}}$ with parameter $\eta$. Letting $\eta = (\eta_x, x\in T)$ denote an element of $H^T$ (parameters in $H$ indexed by elements of the training set), (16.11) can then be written as the maximization of

$$F(\theta,\eta) = \sum_{x\in T}\int_{\mathcal{R}_Z}\log\left(\frac{f_U(x,z;\theta)}{g(z;\eta_x)}\right)g(z;\eta_x)\mu_Z(dz). \tag{16.13}$$

This expression is amenable to a stochastic gradient ascent implementation. We have

$$\partial_\theta\int_{\mathcal{R}_Z}\log\left(\frac{f_U(x,z;\theta)}{g(z;\eta_x)}\right)g(z;\eta_x)\mu_Z(dz) = \int_{\mathcal{R}_Z}\partial_\theta\log f_U(x,z;\theta)g(z;\eta_x)\mu_Z(dz)$$

and

$$\partial_{\eta_x}\int_{\mathcal{R}_Z}\log\left(\frac{f_U(x,z;\theta)}{g(z;\eta_x)}\right)g(z;\eta_x)\mu_Z(dz)$$
$$= \int_{\mathcal{R}_Z}\left(-\partial_\eta\log g(z;\eta_x)g(z;\eta_x) + \log\left(\frac{f_U(x,z;\theta)}{g(z;\eta_x)}\right)\partial_\eta g(z;\eta_x)\right)\mu_Z(dz)$$
$$= \int_{\mathcal{R}_Z}\log\left(\frac{f_U(x,z;\theta)}{g(z;\eta_x)}\right)\partial_\eta\log g(z;\eta_x)g(z;\eta_x)\mu_Z(dz)$$

Here, we have used the fact that, for all $\eta$,

$$\int_{\mathcal{R}_Z}\partial_\eta\log g(z;\eta)g(z;\eta)\,\mu_Z(dz) = \int_{\mathcal{R}_Z}\partial_\eta g(z;\eta)\mu_Z(dz) = 0$$

since $\int_{\mathcal{R}_Z} g(x, \eta) \mu_Z(dz) = 1$.

Denote by $\pi_\eta$ the probability distribution of the random variable $\boldsymbol{Z}$ taking values in $\mathcal{R}_Z^{|T|}$ obtained by sampling $\boldsymbol{Z} = (Z_x, x \in T)$ such that the components $Z_x$ are independent and with p.d.f. $g(\cdot; \eta_x)$ with respect to $\mu_Z$. Define

$$\Phi_1(\theta, \boldsymbol{z}) = \sum_{x \in T} \partial_\theta \log f_U(x, z_x; \theta)$$

and

$$\Phi_2(theta, \boldsymbol{\eta}, \boldsymbol{z}) = \sum_{x \in T} \log\left(\frac{f_U(x, z; \theta)}{g(z; \eta_x)}\right) \partial_\eta \log g(z_x; \eta_x).$$

Then, following section 3.3, one can maximize (16.13) using the algorithm

$$\begin{cases} \theta_{n+1} = \theta_n + \gamma_{n+1} \Phi_1(\theta_n, \boldsymbol{Z}_{n+1}) \\ \boldsymbol{\eta}_{n+1} = \boldsymbol{\eta}_n + \gamma_{n+1} \Phi_2(\theta_n, \boldsymbol{\eta}_n, \boldsymbol{Z}_{n+1}) \end{cases} \tag{16.14}$$

where $\boldsymbol{Z}_{n+1} \sim \pi_{\boldsymbol{\eta}_n}$.

Alternatively (for example when $T$ is large), one can also sample from $x \in T$ at each update. This would require defining $\pi_\eta$ as the distribution on $T \times \mathcal{R}_Z$ with p.d.f. $\varphi_\eta(x, z) = g(z; \eta_x)/N$, where $N = |T|$. One can now use

$$\Phi_1(\theta, x, z) = \partial_\theta \log f_U(x, z; \theta)$$

and

$$\Phi_2(\theta, \eta, z) = \log\left(\frac{f_U(x, z; \theta)}{g(z; \eta)}\right) \partial_\eta \log g(z; \eta),$$

one can use

$$\begin{cases} \theta_{n+1} = \theta_n + \gamma_{n+1} \partial_\theta \log f_U(X_{n+1}, Z_{n+1}; \theta_n) \\ \eta_{n+1, X_{n+1}} = \eta_{n, X_{n+1}} + \gamma_{n+1} \log\left(\frac{f_U(X_{n+1}, Z_{n+1}; \theta_n)}{g(Z_{n+1}; \eta_{n, X_{n+1}})}\right) \partial_\eta \log g(Z_{n+1}; \eta_{n, X_{n+1}}) \end{cases} \tag{16.15}$$

with $(X_{n+1}, Z_{n+1}) \sim \pi_{\boldsymbol{\eta}_n}$. Sampling from a single training sample at each step can be replaced by sampling from a minibatch with obvious modifications.

## 16.5 Remarks

### 16.5.1 Variations on the EM

Based on the formulation of the EM as the solution of (16.6), it should be clear that solving (16.7) at each step can be replaced by any update of the parameter that increases (16.6). For example, (16.7) can be replaced by a partial run of a gradient

ascent algorithm, stopped before convergence. One can also use a coordinate ascent strategy. Assume that $\theta$ can be split into several components, say two, so that $\theta = (\theta^{(1)}, \theta^{(2)})$. Then, (16.7) may then be split into

$$\theta_{n+1}^{(1)} = \underset{\theta^{(1)}}{\mathrm{argmax}} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta^{(1)}, \theta_n^{(2)})\right) f_Z(z \mid x; \theta(n)) \mu_Z(dz)$$

$$\theta_{n+1}^{(2)} = \underset{\theta^{(2)}}{\mathrm{argmax}} \sum_{x \in T} \int_{\mathcal{R}_Z} \log\left(f_U(x, z; \theta_{n+1}^{(1)}, \theta^{(2)})\right) f_Z(z \mid x; \theta(n)) \mu_Z(dz).$$

Doing so is, in particular, useful when both these steps are explicit, but not (16.7).

## 16.5.2   Direct minimization

While the EM algorithm is widely used in the context of partial observations, it is also possible to make explicit the derivative of

$$\log f_X(x; \theta) = \log \int_{\mathcal{R}_Z} f_U(x, z; \theta) \mu_Z(dz)$$

with respect to the parameter $\theta$. Indeed, differentiating the integral and writing $\partial_\theta f_U = f_U \partial_\theta \log f_U$, we have

$$\partial_\theta \log f_X(x; \theta) = \int_{\mathcal{R}_Z} \partial_\theta \log f_U(x, z; \theta) \frac{f_U(x, z; \theta)}{f_X(x; \theta)} \mu_Z(dz)$$

$$= \int_{\mathcal{R}_Z} \partial_\theta \log f_U(x, z; \theta) f_Z(z \mid x, \theta) \mu_Z(dz).$$

In other terms, the derivative of the log-likelihood of the observed data is the conditional expectation of the derivative of the log-likelihood of the full data given the observed data. When computable, this expression can be used with standard gradient-based optimization methods, such as those described in chapter 3. This expression is also amenable to a stochastic gradient ascent algorithm, namely

$$\theta_{n+1} = \theta_n + \gamma_{n+1} \sum_{x \in T} \partial_\theta f_U(x, Z_{n+1,x}, \theta_n) \tag{16.16}$$

where $Z_{n+1,x}$ follows the distribution with density $f_Z(\cdot \mid x, \theta_n)$ with respect to $\mu_Z$. An alternative SGA implementation can use the discussion in section 16.4.4, with the density $g_{\eta_x}$ replaces by $f_Z(\cdot \mid x, \eta_x)$, which leads to

$$\begin{cases} \theta_{n+1} = \theta_n + \gamma_{n+1} \sum_{x \in T} \partial_\theta \log f_U(x, Z_{n+1,x}, \theta_n) \\ \eta_{n+1,x} = \eta_{n,x} - \gamma_{n+1} \partial_{\eta_x} \log f_Z(Z_{n+1,x} \mid x, \eta_x), \quad x \in T \end{cases}$$

where $Z_{n+1,x}$ follows the distribution with density $f_Z(\cdot \mid x, \eta_{n,x})$.

### 16.5.3  Product measure assumption

We have worked, in this chapter, under the assumption that $\pi_U$ was absolutely continuous with respect to a product measure $\mu_U = \mu_X \otimes \mu_Z$. This is not a mild assumption, as it fails to include some important cases, for example when $X$ and $Z$ have some deterministic relationship, the simplest instance being when $X = F(Z)$ for some function $F$. In many cases, however, one can make simple transformations on the model that will make it satisfy this working assumption. For example, if $X = F(Z)$, one can generally split $Z$ into $Z = (Z^{(1)}, Z^{(2)})$ so that the equation $X = F(Z)$ is equivalent to $Z^{(2)} = G(X, Z^{(1)})$ for some function $G$. One can then apply the discussion above to $U = (X, Z^{(1)})$ instead of $U = (X, Z)$.

If one is ready to step further into measure theoretic concepts, however, one can see that this product decomposition assumption was in fact unnecessary. Indeed, one can assume that the measure $\mu_U$ can "disintegrate" in the following sense: there exists, a measure $\mu_X$ on $\mathcal{R}_X$ and, for all $x \in \mathcal{R}_X$, a measure $\mu_Z(\cdot|x)$ on $\mathcal{R}_Z$ such that, for all functions $\psi$ defined on $\mathcal{R}_U$,

$$\int_{\mathcal{R}_U} \psi(x,z)\mu_U(dx,dz) = \int_{\mathcal{R}_X} \int_{\mathcal{R}_Z} \psi(x,z)\mu_Z(dz|x)\mu_X(dx).$$

This is in fact a fairly general situation [34] as soon as one assumes that $\mu_U(\mathcal{R})$ is finite (which is not a real loss of generality as one can reduce to this case by replacing if needed $\mu_U$ by an equivalent probability distribution).

With this assumption, the marginal distribution of $X$ had a p.d.f. with respect to $\mu_X$ given by

$$f_X(x) = \int_{\mathcal{R}_Z} f_U(x,z)\mu_Z(dz|x)$$

and the conditional distributions $P_Z(\cdot \mid x)$ have a p.d.f. relative to $\mu_Z(\cdot \mid x)$ given by

$$f_Z(z \mid x) = \frac{f_U(x,z)}{f_X(x)}.$$

The computations and approximations made earlier in this chapter can then be applied with essentially no modification.

# Chapter 17

# Learning Graphical Models

We discuss, in this chapter, several methods designed to learn parameters of graphical models, starting with the somewhat simpler case of Bayesian networks, than passing to Markov random fields on loopy graphs.

## 17.1 Learning Bayesian networks

### 17.1.1 Learning a Single Probability

Since Bayesian networks are specified by probabilities and conditional probabilities of configurations of variables, we start with a discussion of the basic problem of estimating discrete probability distributions.

The obvious way to estimate the probability of an event $A$ based on a series of $N$ independent experiments is by using relative frequencies

$$f_A = \frac{\#\{A \text{ occurs}\}}{N}.$$

This estimation is unbiased ($\mathbb{E}(f_A) = \mathbb{P}(A)$) and its variance is $\mathbb{P}(A)(1 - \mathbb{P}(A))/N$. This implies that the relative error $\delta_A = f_A/\mathbb{P}(A) - 1$ has zero mean and variance

$$\sigma^2 = \frac{1 - \mathbb{P}(A)}{N\mathbb{P}(A)}.$$

This number can clearly become very large when $\mathbb{P}(A) \simeq 0$. In particular, when $\mathbb{P}(A)$ is small compared to $1/N$, the relative frequency will often be $f_A = 0$, leading to the false conclusion that $A$ is not just rare, but impossible. If there are reasons to expect beforehand that $A$ is indeed possible, it is important to inject this prior belief in the procedure, which suggest using Bayesian estimation methods.

The main assumption for these methods is to consider the unknown probability, $p = \mathbb{P}(A)$, as a random variable, yielding a generative process in which a random probability is first obtained, and then $N$ instances of $A$ or not-$A$ are generated using this probability.

Assume that the "prior distribution" of $p$ (which determines a prior belief) has a p.d.f. $q$ (with respect to Lebesgue's measure) on the unit interval. Given on $N$ independent observations of occurrences of $A$, each following a Bernoulli distribution $b(p)$, the joint likelihood of all involved variables is given by

$$\binom{N}{k} p^k (1-p)^{N-k} q(p),$$

where $k$ is the number of times the event $A$ has been observed.

The conditional density of $p$ given the observation ($k$ occurrences of $A$) is called the *posterior distribution*. Here, it is given by

$$q(p \mid k) = \frac{q(p)}{C_k} p^k (1-p)^{N-k}$$

where $C_k$ is a normalizing constant. If there was no specific prior knowledge on $p$ (so that $q(p) = 1$), the resulting distribution is a beta distribution with parameters $k + 1$ and $N - k + 1$, the beta distribution being defined as follows.

**Definition 17.1** *The beta distribution with parameters $a$ and $b$ (abbreviated $\beta(a, b)$) has density with respect to Lebesgue's measure*

$$\rho(t) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} t^{a-1} (1-t)^{b-1} \text{ if } t \in [0,1]$$

*and $\rho(t) = 0$ otherwise, with*

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt.$$

From the definition of a beta distribution, it is clear also that, if we choose the prior to be $\beta(a + 1, \nu - a + 1)$ then the posterior is $\beta(k + a + 1, N + \nu - (k + a) + 1)$. The posterior therefore belongs to the same family of distributions as the prior, and one says that the beta distribution is a *conjugate prior* for the binomial distribution. The mode of the posterior distribution (which is the maximum a posteriori (MAP) estimator) is given by

$$\hat{p} = \frac{k+a}{N+\nu}.$$

This estimator now provides a positive value even if $k = 0$. By selecting $a$ and $\nu$, one therefore includes the prior belief that $p$ is positive.

### 17.1.2 Learning a Finite Probability Distribution

Now assume that $F$ is a finite space and that we want to estimate a probability distribution $p = (p(x), x \in F)$ using a Bayesian approach as above. We cannot use the previous approach to estimate each $p(x)$ separately, since these probabilities are linked by the fact that they sum to 1. We can however come up with a good (conjugate) prior, identified, as done above, by computing the posterior associated to a uniform prior distribution.

Letting $N_x$ be the number of times $x \in F$ is observed among $N$ independent samples of a random variable $X$ with distribution $P_X(\cdot) = p(\cdot)$, the joint distribution of $(N_x, x \in F)$ is multinomial, given by

$$\mathbb{P}(N_x, x \in F \mid p(\cdot)) = \frac{N!}{\prod_{x \in F} N_x!} \prod_{x \in F} p(x)^{N_x}.$$

The posterior distribution of $p(\cdot)$ given the observations with a uniform prior is proportional to $\prod_{x \in F} p(x)^{N_x}$. It belongs to the family of *Dirichlet distributions*, described in the following definition.

**Definition 17.2** *Let $F$ be a finite set and $\mathcal{S}_F$ be the simplex defined by*

$$\mathcal{S}_F = \left\{ (p(x), x \in F) : p(x) \geq 0, x \in F \text{ and } \sum_{x \in F} p(x) = 1 \right\}.$$

*The Dirichlet distribution with parameters $a = (a(x), x \in F)$ (abbreviated $\mathrm{Dir}(a)$) has density*

$$\rho(p(\cdot)) = \frac{\Gamma(\nu)}{\prod_{x \in F} \Gamma(a(x))} \prod_{x \in F} p(x)^{a(x)-1}, \text{ if } x \in \mathcal{S}_F$$

*and 0 otherwise, with $\nu = \sum_{x \in F} a(x)$.*

Note that, if $F$ has cardinality 2, the Dirichlet distribution coincides with the beta distribution. Similarly to the beta for the binomial, and almost by construction, the Dirichlet distribution is a conjugate prior for the multinomial. More precisely, if the prior distribution for $p(\cdot)$ is $\mathrm{Dir}(1+a(x), x \in F)$, then the posterior after $N$ observations of $X$ is $\mathrm{Dir}(1 + N_x + a(x), x \in F)$, and the MAP estimator is given by

$$\hat{p}(x) = \frac{N_x + a(x)}{N + \nu}$$

with $\nu = \sum_{x \in F} a(x)$.

### 17.1.3   Conjugate Prior for Bayesian Networks

We now consider a Bayesian network on the set $\mathcal{F}(V)$ containing configurations $x = (x^{(s)}, s \in V)$ with $x^{(s)} \in F_s$. We want to estimate the conditional probabilities in the representation

$$\mathbb{P}(X = x) = \prod_{s \in V} p_s(x^{(pa(s))}, x^{(s)}).$$

Assume that $N$ independent observations of $X$ have been made. Define the counts $N_s(x^{(s)}, x^{(pa(s))})$ to be the number of times the observation $x^{(\{s\} \cup pa(s))}$ has been made. Then, it is straightforward to see that, assuming a uniform prior for the $p_s$, their posterior distribution is proportional to

$$\prod_{s \in V} \prod_{x^{(pa(s))} \in F_{pa(s)}} \prod_{x^{(s)} \in F_s} p_s(x^{(pa(s))}, x^{(s)})^{N_s(x^{(s)}, x^{(s^-)})}.$$

This implies that, for the posterior distribution, the conditional probabilities $p_s(x^{(pa(s))}, \cdot)$ are independent and follow a Dirichlet distribution with parameters $1 + N_s(x^{(s)}, x^{(pa(s))})$, $x^{(s)} \in F_s$.

So, independent Dirichlet distributions indexed by configurations of parents of nodes provide a conjugate prior for the general Bayesian network model. This prior is specified by a family of positive numbers

$$\left( a_s(x^{(s)}, x^{(pa(s))}), s \in V, x^{(s)} \in F_s, x^{(pa(s))} \in \mathcal{F}(pa(s)) \right), \tag{17.1}$$

yielding a prior probability proportional to

$$\prod_{s \in V} \prod_{x^{(pa(s))} \in F_{pa(s)}} \prod_{x^{(s)} \in F_s} p_s(x^{(pa(s))}, x^{(s)})^{a_s(x^{(s)}, x^{(s^-)}) - 1}.$$

and a MAP estimator

$$\hat{p}_s(x^{(pa(s))}, x^{(s)}) = \frac{N_s(x^{(s)}, x^{(pa(s))}) + a_s(x^{(s)}, x^{(s^-)})}{N_s(x^{(s^-)}) + \nu_s(x^{(s^-)})} \tag{17.2}$$

where $N_s(x^{(pa(s))}) = \sum_{x^{(s)} \in F_s} N_s(x^{(s)}, x^{(pa(s))})$ and $\nu_s(x^{(pa(s))}) = \sum_{x^{(s)} \in F_s} a_s(x^{(s)}, x^{(pa(s))})$.

One can restrict the huge class of coefficients described by (17.1) to a smaller class by imposing the following condition.

**Definition 17.3** *One says that the family of coefficients*

$$a = (a_s(x^{(s)}, x^{(pa(s))}), s \in V, x^{(s)} \in F_s, x^{(pa(s))} \in \mathcal{F}(pa(s))),$$

*is consistent if there exists a positive scalar $\nu$ and a probability distribution $P'$ on $\mathcal{F}(V)$ such that*

$$a_s(x^{(s)}, x^{(pa(s))}) = \nu P'_{\{s\} \cup pa(s)}(x^{(\{s\} \cup pa(s))}).$$

The class of products of Dirichlet distributions with consistent families of coefficients still provides a conjugate prior for Bayesian networks (the proof being left to the reader). Within this class, the simplest choice (and most natural in the absence of additional information) is to assume that $P'$ is uniform, so that

$$a_s(x^{(s)}, x^{(pa(s))}) = \frac{\nu'}{|\mathcal{F}(\{s\} \cup pa(s))|}. \tag{17.3}$$

With this choice, $\nu'$ is the only parameter that needs to be specified. It is often called the equivalent sample size for the prior distribution.

We can see from (17.2) that using a prior distribution is quite important for Bayesian networks, since, when the number of parents increases, some configurations on $\mathcal{F}(pa(s))$ may not be observed, resulting in an undetermined value for the ratio

$$N_s(x^{(s)}, x^{(pa(s))})/N_s(x^{(s^-)}),$$

even though, for the estimated model, the probability of observing $x^{(pa(s))}$ may not be zero.

### 17.1.4 Structure Scoring

Given a prior defined as a family of Dirichlet distributions associated to $a = (a_s(x^{(s)}, x^{(pa(s))})$ for $s \in V, x^{(s)} \in F_s, x^{(pa(s))} \in \mathcal{F}(pa(s))$, the joint density of the observations and parameters is given by

$$P(x, \theta) = \prod_{s, x^{(pa(s))}} \mathcal{D}(a_s(\cdot, x^{(pa(s))})) \prod_{s, x^{(s)}, x^{(pa(s))}} p(x^{(pa(s))}, x^{(s)})^{N_s(x^{(s)}, x^{(pa(s))}) + a_s(x^{(s)}, x^{(pa(s))}) - 1}$$

with

$$\mathcal{D}(a(\lambda), \lambda \in F) = \frac{\Gamma(\nu)}{\prod_\lambda \Gamma(a(\lambda))}$$

and $\nu = \sum_\lambda a(\lambda)$. Here, $\theta$ represents the parameters of the model, i.e., the conditional distributions that specify the Bayesian network. Note that $P(x, \theta)$ is a density over the product space $\mathcal{F}(V) \times \Theta$ where $\Theta$ is the space of all these conditional distributions. The marginal of this likelihood over all possible parameters, i.e.,

$$P(x) = \int P(\mathbf{x}, \theta) d\theta$$

provides the expected likelihood of the sample relative to the distribution of the parameters, and only depends on the structure of the network. In our case, integrating with respect to $\theta$ yields

$$\log P(x) = \sum_{s, x_{pa(s)}} \log \frac{\mathcal{D}(a_s(\cdot, x^{(pa(s))}))}{\mathcal{D}(a_s(\cdot, x^{(pa(s))}) + N_s(\cdot, x^{(pa(s))}))}.$$

Letting

$$\gamma(s, pa(s)) = \sum_{x^{(pa(s))}} \log \frac{\mathcal{D}(a_s(\cdot, x^{(pa(s))}))}{\mathcal{D}(a_s(\cdot, x^{(pa(s))}) + N_s(\cdot, x^{(pa(s))}))},$$

the decomposition

$$\log P(x) = \sum_{s \in V} \gamma(s, pa(s))$$

expresses this likelihood as a sum of "scores" (associated to each node and its parents), which depends on the observed sample. The scores that are computed above are often called *Bayesian scores* because they derive from a Bayesian construction. One can also consider simpler scores, such as penalized likelihood:

$$\gamma(s, pa(s)) = - \sum_{x^{(pa(s))}} \hat{\mathcal{H}}(X^{(s)} \mid X^{(pa(s))}) |\mathcal{F}(pa(s))| - \rho |pa(s)|,$$

where $\hat{\mathcal{H}}$ is the conditional entropy for the empirical distribution based on observed samples. Structure learning algorithms [144, 108] are designed to optimize such scores.

### 17.1.5   Reducing the Parametric Dimension

In the previous section, we estimated all conditional probabilities intervening in the network. This is obviously a lot of parameters and, even with a regularizing prior, the estimated values are likely to be be inaccurate for small sample sizes. It then becomes desirable to simplify the parametric complexity of the model.

When the sets $F_s$ are not too large, which is common in practice, the parametric explosion is due to the multiplicity of parents, since the number of conditional probabilities $p_s(x^{(pa(s))}, \cdot)$ grows exponentially with $|pa(s)|$. One way to simplify this is to assume that the conditional probability at $s$ only depends on $x^{(pa(s))}$ via some "global-effect" statistic $g_s(x^{(pa(s))})$. The idea, of course, is that the number of values taken by $g_s$ should remain small, even if the number of parents is large.

Examples of some functions $g_s$ can be $\max(x^{(t)}, t \in pa(s))$, or the min, or some simple (quantized) function of the sum. With binary variables ($F_s = \{0, 1\}$), logical operators are also available ("and", "or", "xor"), as well as combinations of them. The choice made for the functions $g_s$ is part of building the model, and would rely on the specific context and prior information on the process, which is always important to account for, in any statistical problem.

Once the $g_s$'s are fixed, learning the network distribution, which is now given by

$$\pi(x) = \prod_{s \in V} p_s(g_s(x^{(pa(s))}), x^{(s)})$$

can be done exactly as before, the parameters being all $p_s(w, \lambda), \lambda \in F_s, w \in W_s$, where $W_s$ is the range of $g_s$, and Dirichlet priors can be associated to each $p_s(w, \cdot)$ for $s \in V$ and $w \in W_s$. The counts provided in (17.3) now can be chosen as

$$a_s(x_s, w) = \frac{\nu'}{|F||g_s^{-1}(w)|}. \tag{17.4}$$

## 17.2 Learning Loopy Markov Random Fields

Like everything else, parameter estimation for loopy networks is much harder than with trees or Bayesian networks. There is usually no closed form expression for the estimators, and their computation relies on more or less tractable numerical procedures.

### 17.2.1 Maximum Likelihood with Exponential Models

In this section, we consider a parametrized model for a Gibbs distribution

$$\pi_\theta(x) = \frac{1}{Z_\theta} \exp(-\theta^T U(x)) \tag{17.5}$$

where $\theta$ is a $d$-dimensional parameter and $U$ is a function from $\mathcal{F}(V)$ to $\mathbb{R}^d$. For example, if $\pi$ is an Ising model with

$$\pi(x) = \frac{1}{Z} \exp\left(\alpha \sum_{s \in V} x^{(s)} + \beta \sum_{s \sim t} x^{(s)} x^{(t)}\right),$$

then $\theta = (\alpha, \beta)$ and $U(x) = -(\sum_s x^{(s)}, \sum_{s \sim t} x^{(s)} x^{(t)})$. Most of the Markov random fields models that are used in practice can be put in this form. The constant $Z_\theta$ in (17.5) is

$$Z_\theta = \sum_{x \in \mathcal{F}(V)} \exp(-\theta^T U(x))$$

and is usually not computable.

Now, assume that an $N$-sample, $x_1, \ldots, x_N$, is observed for this distribution. The maximum likelihood estimator maximizes

$$\ell(\theta) = \frac{1}{N} \sum_{k=1}^{N} \log \pi_\theta(x_k) = -\theta^T \bar{U}_N - \log Z_\theta$$

with $\bar{U}_N = (U(x_1) + \cdots + U(x_N))/N$.

We have the following proposition, which is a well-known property of exponential families of probabilities.

**Proposition 17.4** *The log-likelihood, $\ell$, is a concave function of $\theta$, with*

$$\nabla \ell(\theta) = E_\theta(U) - \bar{U}_N \tag{17.6}$$

*and*

$$\nabla^2 \ell(\theta) = -Var_\theta(U) \tag{17.7}$$

*where $E_\theta$ denotes the expectation with respect to $\pi_\theta$ and $Var_\theta$ the covariance matrix under the same distribution.*


We skip the proof, which is just computation. This proposition implies that a local maximum of $\theta \mapsto \ell(\theta)$ must also be global. Any such maximum must be a solution of

$$E_\theta(U) = \bar{U}_N(x_0)$$

and conversely. There are some situations in which the maximum does not exist, or is not unique. Let us first discuss the second case.

If several solutions exist, the log-likelihood cannot be strictly concave: there must exist at least one $\theta$ for which $Var_\theta(U)$ is not definite. This implies that there exists a nonzero vector $u$ such that $var_\theta(u^T U) = u^T Var_\theta(U)u = 0$. This is only possible when $u^T U(x) = \text{cst}$ for all $x \in F_V$. Conversely, if this is true, $Var_\theta(U)$ is degenerate *for all $\theta$*.

So, the non-uniqueness of the solutions is only possible when a deterministic affine relation exists between the components of $U$, i.e., when the model is over-dimensioned. Such situations are usually easily dealt with by removing some parameters. In all other cases, there exists at most one maximum.

For a concave function like $\ell$ to have no maximum, there must exist what is called a direction of recession [167], which is a direction $\alpha \in \mathbb{R}^d$ such that, for all $\theta$, the function $t \mapsto \ell(\theta + t\alpha)$ is increasing. In this case the maximum is attained "at infinity". Denoting $U_\alpha(x) = \alpha^T U(x)$, the derivative in $t$ of $\ell(\theta + t\alpha)$ is

$$E_{\theta + t\alpha}(U_\alpha) - \bar{U}_\alpha$$

where $\bar{U}_\alpha = \alpha^T \bar{U}_N$. This derivative is positive for all $t$ if and only if

$$\bar{U}_\alpha = U_\alpha^* := \min\{U_\alpha(x), x \in \mathcal{F}(V)\} \tag{17.8}$$

and $U_\alpha$ is not constant. To prove this, assume that the derivative is positive. Then $U_\alpha$ is not constant (otherwise, the derivative would be zero). Let $\mathcal{F}_\alpha^* \subset \mathcal{F}(V)$ be the

set of configurations $x$ for which $U_\alpha(x) = U_\alpha^*$. Then

$$
\begin{aligned}
&E_{\theta+t\alpha}(U_\alpha) \\
&= \frac{\sum_{x\in\mathcal{F}(V)} U_\alpha(x)\exp(-\theta^T U(x) - tU_\alpha(x))}{\sum_{x\in\mathcal{F}(V)}\exp(-\theta^T U(x) - tU_\alpha(x))} \\
&= \frac{\sum_{x\in\mathcal{F}(V)} U_\alpha(x)\exp(-\theta^T U(x) - t(U_\alpha(x) - U_\alpha^*))}{\sum_{x\in\mathcal{F}(V)}\exp(-\theta^T U(x) - t(U_\alpha(x) - U_\alpha^*))} \\
&= \frac{U_\alpha^*\sum_{x\in\mathcal{F}_\alpha^*}\exp(-\theta^T U(x)) + \sum_{x\notin\mathcal{F}_\alpha^*} U_\alpha(x)\exp(-\theta^T U(x) - t(U_\alpha(x) - U_\alpha^*))}{\sum_{x\in\mathcal{F}_\alpha^*}\exp(-\theta^T U(x)) + \sum_{x\notin\mathcal{F}_\alpha^*}\exp(-\theta^T U(x) - t(U_\alpha(x) - U_\alpha^*))}.
\end{aligned}
$$

When $t$ tends to $+\infty$, the sums over $x \notin \mathcal{F}_\alpha^*$ tend to 0, which implies that $E_{\theta+t\alpha}(U_\alpha)$ tends to $U_\alpha^*$. So, if $E_{\theta+t\alpha}(U_\alpha) - \bar{U}_\alpha > 0$ for all $t$, then $\bar{U}_\alpha = U_\alpha^*$ and $U_\alpha$ is not constant. The converse statement is obvious.

As a conclusion, the function $\ell$ has a finite maximum if and only if there is no direction $\alpha \in \mathbb{R}^d$ such that $\alpha^T(U(x) - \bar{U}_N) \le 0$ for all $x \in \mathcal{F}(V)$. Equivalently, $\bar{U}_N$ must belong to the interior of the convex hull of the finite set

$$
\{U(x), x \in \mathcal{F}(V)\} \subset \mathbb{R}^d.
$$

In such a case, that we hereafter assume, computing the maximum likelihood estimator boils down to solving the equation

$$
E_\theta(U) = \bar{U}_N.
$$

Because the maximization problem is concave, we know that numerical algorithms such as gradient ascent,

$$
\theta(t+1) = \theta(t) + \epsilon(E_{\theta(t)}(U) - \bar{U}_N), \tag{17.9}
$$

converge to the optimal parameter. Unfortunately, the computation of the expectations and covariance matrices can only be made explicitly for acyclic models, for which parameter estimation is not a problem anyway. For general loopy graphical models, the expectation can be estimated iteratively using Monte-Carlo methods. It turns out that this estimation can be synchronized with gradient descent to obtain a consistent algorithm, which is described in the next section.

### 17.2.2 Maximum likelihood with stochastic gradient ascent

As remarked above, for fixed $\theta$, we have designed, in chapter 14, Markov chain Monte Carlo algorithms that asymptotically sample form $\pi_\theta$. Select one of these algorithms, and let $p_\theta$ be the corresponding transition probabilities for a given $\theta$, so

that $p_\theta(x, y) = \mathbb{P}(X_{n+1} = y \mid X_n = x)$ for the sampling chain. Then, define the iterative algorithm, initialized with arbitrary $\theta_0$ and $x_0 \in \mathcal{F}(V)$, that loops over the following two steps.

(SG1)  Sample from the distribution $p_{\theta_t}(x_t, \cdot)$ to obtain a new configuration $x_{t+1}$.

(SG2)  Update the parameter using

$$\theta_{t+1} = \theta_t + \gamma_{t+1}(U(x_{t+1}) - \bar{U}_N). \tag{17.10}$$

This algorithm differs from the situation considered in section 3.3 in that the distribution of the sampled variable $x_{t+1}$ depends on both the current parameter $\theta_t$ and on the current variable $x_t$. Convergence requires additional constraints on the size of the gains $\gamma(t)$ and we have the following theorem [206].

**Theorem 17.5**  *If $p_\theta$ corresponds to the Gibbs sampler or Metropolis algorithm, and $\gamma_{t+1} = \epsilon/(t+1)$ for small enough $\epsilon$, the algorithm that iterates (SG1) and (SG2) converges almost surely to the maximum likelihood estimator.*

The speed of convergence of such algorithms depends both on the speed of convergence of the Monte-Carlo sampling and of the original gradient ascent. The latter can be improved somewhat with variants similar to those discussed in section 3.3, for example by choosing data-adaptive gains as in the ADAM algorithm.

### 17.2.3   Relation with Maximum Entropy

The maximum likelihood estimator is closely related to what is called the maximum entropy extension of a set of constraints. Let the function $U$ from $\mathcal{F}(V)$ to $\mathbb{R}^d$ be given. An element $u \in \mathbb{R}^d$ is said to be a consistent assignment for $U$ if there exists a probability distribution $\pi$ on $\mathcal{F}(V)$ such that $E_\pi(U) = u$. An example of consistent assignment is any empirical average $\bar{U}$ based on a sample $(x^{(1)}, \ldots, x^{(N)})$, since $\bar{U} = E_\pi(U)$ for

$$\pi = \frac{1}{N} \sum_{k=1}^{N} \delta_{x^{(k)}}.$$

Given $U$ and a consistent assignment, $u$, the associated maximum entropy extension is defined as a probability distribution $\pi$ maximizing the entropy, $\mathcal{H}(\pi)$, subject to the constraint $E_\pi(U) = u$. This is a convex optimization problem, with constraints

$$\begin{cases} \displaystyle\sum_{x \in \mathcal{F}(V)} \pi(x) = 1 \\ \displaystyle\sum_{x \in \mathcal{F}(V)} U_j(x)\pi(x) = u_j, j = 1, \ldots, d \\ \pi(x) \geq 0, x \in \mathcal{F}(V) \end{cases} \tag{17.11}$$

Because the entropy is strictly convex, there is a unique solution to this problem. We first discuss non-positive solutions, i.e., solutions for which $\pi(x) = 0$ for some $x$. An important fact is that, if, for a given $x$, there exists $\pi_1$ such that $E_{\pi_1}(U) = u$ and $\pi_1(x) > 0$, then the optimal $\pi$ must also satisfy $\pi(x) > 0$. This is because, if $\pi(x) = 0$, then, letting $\pi_\epsilon = (1 - \epsilon)\pi + \epsilon\pi_1$, we have $E_{\pi_\epsilon}(U) = u$ since this constraint is linear, $\pi_\epsilon(x) > 0$ and

$$
\begin{aligned}
H(\pi_\epsilon) - H(\pi) &= -\sum_{y,\pi(y)>0} (\pi_\epsilon(y)\log\pi_\epsilon(y) - \pi(y)\log\pi(y)) \\
&\quad - \sum_{y,\pi(y)=0} \epsilon\pi_1(y)(\log(\epsilon) + \log\pi_1(y)) \\
&= -\epsilon\log\epsilon \sum_{y,\pi(y)=0} \pi_1(y) + O(\epsilon)
\end{aligned}
$$

which is positive for small enough $\epsilon$, contradicting the fact that $\pi$ is a maximizer.

Introduce the set $\mathcal{N}_u$ containing all configurations $x \in \mathcal{F}(V)$ such that $\pi(x) = 0$ for all $\pi$ such that $E_\pi(U) = u$. Then we know that the maximum entropy extension satisfies $\pi(x) > 0$ if $x \notin \mathcal{N}_u$. Introduce Lagrange multipliers $\theta_0, \theta_1, \ldots, \theta_d$ for the $d + 1$ equality constraints in (17.11), and the Lagrangian

$$
L = H(\pi) + \sum_{x \in \mathcal{F}(V)\setminus\mathcal{N}_u} (\theta_0 + \theta^T U(x))\pi(x)
$$

in which we have set $\theta = (\theta_1, \ldots, \theta_d)$, we find that the optimal $\pi$ must satisfy

$$
\begin{cases}
\log\pi(x) = -\theta_0 - 1 - \theta^T U(x) \\
\sum_x \pi(x) = 1 \\
E_\pi(U) = \bar{u}
\end{cases}
$$

In other terms, the maximum entropy extension is characterized by

$$
\pi(x) = \frac{1}{Z_\theta} \exp(-\theta^T U(x))\mathbf{1}_{\mathcal{N}_u^c}(x)
$$

and $E_\pi(U) = u$.

In particular, if $\mathcal{N}_u = \emptyset$, then the maximum entropy extension is positive. If, in addition, $u = \bar{U}$ for some observed sample, then it coincides with the maximum likelihood estimator for (17.5). Notice that, in this case, the condition $\mathcal{N}_u \neq \emptyset$ coincide with the condition that there exists $\alpha$ such that $\alpha^T U(x) \geq \alpha^T u$ for all $x$, with $\alpha^T U(x)$ not constant. Indeed, assume that the latter condition is true. Then, if $E_\pi(U) = u$, then $E_\pi(\alpha^T U) = \alpha^T u$, which is only possible if $\pi(x) = 0$ for all $x$ such

that $\alpha^T U(x) < \alpha^T u$. Such $x$'s exist by assumption, and therefore $\mathcal{N}_u \neq \emptyset$. Conversely, assume $\mathcal{N}_u \neq \emptyset$. If condition (17.8) is not satisfied, then we have shown when discussing maximum likelihood that an optimal parameter for the exponential model would exist, leading to a positive distribution for which $E_\pi(U) = u$, which is a contradiction.

### 17.2.4  Iterative Scaling

Iterative scaling is a method that is well-adapted to learning distributions given by (17.5), when $U$ can be interpreted as a random histogram, or a collection of them. More precisely, assume that for all $x \in \mathcal{F}(V)$, one has

$$U(x) = (U_1(x), \dots, U_q(x))$$

with

$$\sum_{j=1}^{q} U_j(x) = 1 \text{ and } U_j(x) \geq 0.$$

Let the parameter be given by $\theta = (\theta_1, \dots, \theta_q)$. Assume that $x_1, \dots, x_N$ have been observed, and let $u \in \mathbb{R}^d$ be a consistent assignment for $U$, with $u_j > 0$ for $j = 1, \dots, d$ and such that $\mathcal{N}_u = \emptyset$. Iterative scaling computes the maximum entropy extension of $E_\pi(U) = u$, that we will denote $\pi^*$. It is supported by the following lemma.

**Lemma 17.6** *Let $\pi$ be a probability on $\mathcal{F}(V)$ with $\pi > 0$ and define*

$$\pi'(x) = \frac{\pi(x)}{\zeta} \prod_{j=1}^{d} \left( \frac{u_j}{E_\pi(U_j)} \right)^{U_j(x)}$$

*where $\zeta$ is chosen so that $\pi'$ is a probability. Then $\pi' > 0$ and*

$$KL(\pi^* \| \pi') - KL(\pi^* \| \pi) \leq -KL(u \| E_\pi(U)) \leq 0 \tag{17.12}$$

Proof  Note that, since $\pi > 0$, $E_\pi(U_j)$ must also be positive for all $j$, since $E_\pi(U_j) = 0$ would otherwise imply $U_j = 0$ and $u_j = 0$ for $u$ to be consistent. So, $\pi'$ is well defined and obviously positive.

We have

$$
\begin{aligned}
KL(\pi^* \| \pi') - KL(\pi^* \| \pi) &= \log \zeta - \sum_{x \in \mathcal{F}(V)} \pi^*(x) \sum_{j=1}^{d} U_j(x) \log \frac{u_j}{E_\pi(U_j)} \\
&= \log \zeta - \sum_{j=1}^{d} u_j \log \frac{u_j}{E_\pi(U_j)} \\
&= \log \zeta - KL(u \| E_\pi(U)).
\end{aligned}
$$

(We have used the identity $E_{\pi^*}(U) = u$.) So it suffices to prove that $\zeta \leq 1$. We have

$$
\begin{aligned}
\zeta &= \sum_{x \in \mathcal{F}(V)} \pi(x) \prod_{j=1}^{d} \left( \frac{u_j}{E_\pi(U_j)} \right)^{U_j(x)} \\
&\leq \sum_{j=1}^{d} \sum_{x \in \mathcal{F}(V)} \pi(x) U_j(x) \frac{u_j}{E_\pi(U_j)} \\
&= \sum_{j=1}^{d} E_\pi(U_j) \frac{u_j}{E_\pi(U_j)} = 1,
\end{aligned}
$$

which proves the lemma (we have used the fact that, for $x_i$, $w_i$ positive numbers with $\sum_i w_i = 1$, one has $\prod_i x_i^{w_i} \leq \sum_i w_i x_i$, which is a consequence of the concavity of the logarithm). ∎

Consider the iterative algorithm

$$
\pi_{n+1}(x) = \frac{\pi_n(x)}{\zeta_n} \prod_{j=1}^{d} \left( \frac{u_j}{E_{\pi_n}(U_j)} \right)^{U_j(x)}
$$

initialized with a uniform distribution. Equivalently, using the exponential formulation, define, for $j = 1, \ldots, d$,

$$
\theta_{n+1,j} = \theta_{n,j} + \log \frac{E_{\theta_n}(U_j)}{u_j} + KL(u \| E_{\theta_n}(U)), \tag{17.13}
$$

with $\pi_\theta$ given by (17.5), initialized with $\theta_0 = 0$. Note that adding a term that is independent of $j$ to $\theta$ does not change the value of $\pi_\theta$, because the $U_j$'s sum to 1. The model is in fact overparametrized, and the addition of the KL divergence in (17.13) ensures that $\sum_{i=1}^{d} u_i \theta_i = 0$ at all steps.

This algorithm always reduces the Kullback-Leibler distance to the maximum entropy extension. This distance being always positive, it therefore converges to a limit, which, still according to lemma 17.6, is only possible if $KL(u \| E_{\pi_n}(U))$ also tends to 0, that is $E_{\pi_n}(U) \to u$. Since the space of probability distributions is compact, the Heine-Borel theorem implies that the sequence $\pi_{\theta_n}$ has at least one accumulation point, that we now identify. If $\pi$ is such a point, one must have $E_\pi(U) = u$. Moreover, we have $\pi > 0$, since otherwise $KL(\pi^* \| \pi) = +\infty$. To prove that $\pi = \pi^*$ (and therefore the limit of the sequence), it remains to show that it can be put in the form (17.5). For this, define the vector space $\mathcal{V}$ of functions $v : \mathcal{F}(V) \to \mathbb{R}$ which can be written in the form

$$
v(x) = \alpha_0 + \sum_{j=1}^{g} \alpha_j U_j(x).
$$

Since $\log \pi_{\theta_n} \in \mathcal{V}$ for all $n$, so is its limit, and this proves that $\log \pi$ belongs to $\mathcal{V}$. We have obtained the following proposition.

**Proposition 17.7** *Assume that for all $x \in \mathcal{F}(V)$, one has $U(x) = (U_1(x),\ldots,U_d(x))$ with*

$$\sum_{j=1}^{d} U_j(x) = 1 \text{ and } U_j(x) \geq 0.$$

*Let $u$ be a consistent assignment for the expectation of $U$ such that $\mathcal{N}_u = \emptyset$. Then, the algorithm described in (17.13) converges to the maximum entropy extension of $u$.*

This is the iterative scaling algorithm. This method can be extended in a straightforward way to handle the maximum entropy extension for a family of functions $U^{(1)},\ldots,U^{(K)}$, such that, for all $x$ and for all $k$, $U^{(k)}(x)$ is a $d_k$-dimensional vector such that

$$\sum_{j=1}^{d_k} U_j^{(k)}(x) = 1.$$

The maximum entropy extension takes the form

$$\pi_\theta(x) = \frac{1}{Z_\theta} \exp\left(-\sum_{k=1}^{K} (\theta^{(k)})^T U^{(k)}(x)\right),$$

where $\theta^{(k)}$ is $d_k$-dimensional, and iterative scaling can then be implemented by updating only one of these vectors at a time, using (17.13) with $U = U^{(k)}$.

The restriction to $U(x)$ providing a discrete probability distribution for all $x$ is, in fact, no loss of generality. This is because adding a constant to $U$ does not change the resulting exponential model in (17.5), and multiplying $U$ by a constant can be also compensated by dividing $\theta$ by the same constant in the same model. So, if $u_-$ is a lower bound for $\min_{j,x} U_j(x)$, one can replace $U$ by $(U - u_-)$, and therefore assume that $U \geq 0$, and if $u_+$ is an upper bound for $\sum_j U_j(x)$, we can replace $U$ by $U/u_+$ and therefore assume that $\sum_j U_j(x) \leq 1$. Define

$$U_{d+1}(x) = 1 - \sum_{j=1}^{d} U_j(x) \geq 0.$$

Then, the maximum entropy extension for $(U_1,\ldots,U_d)$ with assignment $(u_1,\ldots,u_d)$ is obviously also the extension for $(U_1,\ldots,U_{d+1})$, with assignment $(u_1,\ldots,u_{d+1})$, where

$$u_{d+1} = 1 - \sum_{j=1}^{d} u_j,$$

and the latter is in the form required in proposition 17.7. Note that iterative scaling requires to compute the expectation of $U_1, \ldots, U_d$ before each update. These are not necessarily available in closed form and may have to be estimated using Monte-Carlo sampling.

### 17.2.5  Pseudo likelihood

Maximum likelihood estimation is a special case of *minimal contrast estimators*. These estimators are based on the definition of a measure of dissimilarity, say $C(\pi \| \tilde{\pi})$, between two probability distributions $\pi$ and $\tilde{\pi}$. The usual assumptions on $C$ are that $C(\pi \| \tilde{\pi}) \geq 0$, with equality if and only if $\pi = \tilde{\pi}$, and that $C$ is — at least — continuous in $\pi$ and $\tilde{\pi}$. Minimal contrast estimators approximate the problem of minimizing $\theta \mapsto C(\pi_{\text{true}} \| \pi_\theta)$ over a parameter $\theta \in \Theta$, (which is not feasible, since $\pi_{\text{true}}$, the true distribution of the data, is unknown) by the minimization of $\theta \mapsto C(\hat{\pi} \| \pi_\theta)$ where $\hat{\pi}$ is the empirical distribution computed from observed data. Under mild conditions on $C$, these estimators are generally consistent when $N$ tends to infinity, which means that the estimated parameter asymptotically (in the sample size $N$) provides the best (according to $C$) approximation of $\pi_{\text{true}}$ by the family $\pi_\theta, \theta \in \Theta$.

The contrast that is associated with maximum likelihood is the Kullback-Leibler divergence. Indeed, given a sample $x_1, \ldots, x_N$, we have

$$
\begin{aligned}
KL(\hat{\pi} \| \pi_\theta) &= E_{\hat{\pi}} \log \hat{\pi} - E_{\hat{\pi}} \log \pi_\theta \\
&= E_{\hat{\pi}} \log \hat{\pi} - \sum_{k=1}^{N} \log \pi_\theta(x_k).
\end{aligned}
$$

Since $E_{\hat{\pi}} \log \hat{\pi}$ does not depend on $\theta$, minimizing $KL(\hat{\pi} \| \pi_\theta)$ is equivalent to maximizing $\sum_{k=1}^{N} \log \pi_\theta(x_k)$ which is the log-likelihood.

Maximum pseudo-likelihood estimators form another class of minimal contrast estimators for graphical models. Given a distribution $\pi$ on $\mathcal{F}(V)$, define the local specifications $\pi_s(x^{(s)} \mid x^{(t)}, t \neq s)$ to be conditional distributions at one vertex given the others, and the contrast

$$
C(\pi \| \tilde{\pi}) = \sum_{s \in V} E_\pi (\log \frac{\pi_s}{\tilde{\pi}_s}).
$$

Because we can write, using standard properties of conditional expectations,

$$
C(\pi \| \tilde{\pi}) = \sum_{s \in V} E_\pi \left( E_{\pi_s} (\log \frac{\pi_s}{\tilde{\pi}_s}) \right) = \sum_{s \in V} E(KL(\pi_s(\cdot \mid X^{(t)}, t \neq s) \| \tilde{\pi}_s(\cdot \mid X^{(t)}, t \neq s)),
$$

we see that $C(\pi, \tilde{\pi})$ is always positive, and vanishes (under the assumption of positive $\pi$) only if all the local specifications for $\pi$ and $\tilde{\pi}$ coincide, and this can be shown

to imply that $\pi = \tilde{\pi}$. Indeed, for any $x, y \in \mathcal{F}(V)$, and choosing some order $V = \{s_1, \ldots, s_n\}$ on $V$, one can write

$$\frac{\pi(x)}{\pi(y)} = \prod_{k=1}^{n} \frac{\pi(x^{(s_k)}|x^{(s_1)}, \ldots, x^{(s_{k-1})}, y^{(s_{k+1})}, \ldots, y^{(s_n)})}{\pi(x^{(s_k)}|x^{(s_1)}, \ldots, x^{(s_{k-1})}, y^{(s_{k+1})}, \ldots, y^{(s_n)})}$$

and the ratios $\pi(x)/\pi(y)$, for $x \in \mathcal{F}(V)$, combined with the constraint that $\sum_x \pi(x) = 1$ uniquely define $\pi$.

So $C$ is a valid contrast and

$$C(\hat{\pi}\|\pi_\theta) = \sum_{s \in V} E_{\hat{\pi}} \log \hat{\pi}_s - \sum_{s \in V} \sum_{k=1}^{N} \log \pi_{\theta,s}(x_k^{(s)}|x_k^{(t)}, t \neq s).$$

This yields the *maximum pseudo-likelihood estimator* (or pseudo maximum likelihood) defined as a maximizer of the function (called log-pseudo-likelihood)

$$\theta \mapsto \sum_{s \in V} \sum_{k=1}^{N} \log \pi_{\theta,s}(x_k^{(s)}|x_k^{(s)}, t \neq s).$$

Although maximum likelihood is known to provide the most accurate approximations in many cases, maximum of pseudo likelihood has the important advantage to be, most of the time, computationally feasible. This is because, for a model like (17.5), local specifications are given by

$$\pi_{\theta,s}(x^{(s)} \mid x^{(t)}, t \neq s) = \frac{\exp(-\theta^T U(x))}{\sum_{y^{(s)} \in F_s} \exp(-\theta^T U(y^{(s)} \wedge x^{(V \setminus s)}))}.$$

and therefore include no intractable normalizing constant. Maximum of pseudo-likelihood estimators can be computed using standard maximization algorithms. For exponential models such as (17.5), the log-pseudo-likelihood is, like the log-likelihood, a concave function.

### 17.2.6   Continuous variables and score matching

The methods that were presented so far for discrete variables formally generalize to more general state spaces, even though consistency or convergence issues in non-compact cases can be significantly harder to address. Score matching is a parameter estimation method that was introduced in [95] and was designed, in its original version, to estimate parameters for statistical models taking the form

$$\pi_\theta(x) = \frac{1}{C(\theta)} \exp(-F(x, \theta))$$

with $x \in \mathbb{R}^d$. We assume below suitable integrability and differentiability conditions, in order to justify differentiation under integrals whenever they are needed. The "score function" is defined as

$$s(x, \theta) = -\nabla_x \log \pi_\theta(x) = \nabla_x F(x, \theta)$$

where $\nabla_x$ denotes the gradient with respect to the $x$ variable. Letting $\pi_{\text{true}}$ denote the p.d.f. of the true data distribution (not necessarily part of the statistical model), score matching minimizes

$$f(\theta) = \int_{\mathbb{R}^d} |s(x, \theta) - s_{\text{true}}(x)|^2 \pi_{\text{true}}(x) dx$$

where $s_{\text{true}} = -\nabla \log \pi_{\text{true}}$. This integral can be restricted to the support of $\pi_{\text{true}}$, if we don't want to assume that $\pi_{\text{true}}$ is non-vanishing. Note, however that $f(\theta) = 0$ implies that $\log \pi_\theta(\cdot, \theta) = \log \pi_{\text{true}} \ \pi_{\text{true}}$-almost everywhere, so that $\pi_\theta(x) = c\pi_{\text{true}}(x)$ for some constant $c$ and $x$ in the support of $\pi_{\text{true}}$. Only if $\pi_{\text{true}}(x) > 0$ for all $x \in \mathbb{R}^d$, can we conclude that this requires $\pi_\theta = \pi_{\text{true}}$.

Expanding the squared norm and applying the divergence theorem yield

$$f(\theta) = \int_{\mathbb{R}^d} |\nabla_x \log \pi_\theta(x)|^2 \pi_{\text{true}}(x) dx - 2 \int_{\mathbb{R}^d} \nabla_x \log \pi_\theta(x)^T \nabla \pi_{\text{true}}(x) dx$$
$$+ \int_{\mathbb{R}^d} |s_{\text{true}}(x)|^2 \pi_{\text{true}}(x) dx$$
$$= \int_{\mathbb{R}^d} |\nabla_x \log \pi_\theta(x)|^2 \pi_{\text{true}}(x) dx + 2 \int_{\mathbb{R}^d} \Delta \log \pi_\theta(x)^T \pi_{\text{true}}(x) dx + \int_{\mathbb{R}^d} |s_{\text{true}}(x)|^2 dx$$

To justify the use of the divergence theorem, one needs to assume two derivatives in the log-likelihoods with sufficient decay at infinity (see Hyvärinen and Dayan [95] for details). This shows that minimizing $f$ is equivalent to minimizing

$$g(\theta) = \int_{\mathbb{R}^d} |\nabla_x \log \pi_\theta(x)|^2 \pi_{\text{true}}(x) dx + 2 \int_{\mathbb{R}^d} \Delta \log \pi_\theta(x)^T \pi_{\text{true}}(x) dx$$
$$= \mathbb{E}(|\nabla_x \log \pi_\theta(X)|^2 + 2\Delta \log \pi_\theta(X)).$$

In this form, the objective function can be approximated by a sample average, so that, given observed data $x_1, \ldots, x_N$, one can define the score-matching estimator as a minimizer of

$$\sum_{k=1}^{N} \left( |\nabla_x \log \pi_\theta(x_k)|^2 + 2\Delta \log \pi_\theta(x_k) \right). \tag{17.14}$$

**Remark 17.8** The method can be adapted to deal with discrete variables replacing derivatives with differences. Let $X$ take values in a finite set, $\mathcal{R}_X$, on which a graph

structure can be defined, writing $x \sim y$ if $x$ and $y$ are connected by an edge. For example, if $X$ is itself a Markov random field on a graph $\mathcal{G} = (V, E)$, so that $\mathcal{R}_X = \mathcal{F}(V)$, one can define $x \sim y$ if and only if $x^{(s)} = y^{(s)}$ for all but one $s \in V$. One can then define the score function

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

defined over all $x, y \in \mathcal{R}_X$ such that $x \sim y$. Now the score matching functional is

$$f(\theta) = \sum_{x \in \mathcal{R}_X} \left( \sum_{y \sim x} |s_\theta(x, y) - s_*(x, y)|^2 \right) \pi_*(x),$$

whose minimization is, after reordering terms, equivalent to that of

$$g(\theta) = \sum_{x \in \mathcal{R}_X} \sum_{y \sim x} \left| 1 - \frac{\pi_\theta(y)}{\pi_\theta(x)} \right|^2 \pi_*(x) + 2 \sum_{x \in \mathcal{R}_X} \sum_{y \sim x} \left( \frac{\pi_\theta(x)}{\pi_\theta(y)} - \frac{\pi_\theta(y)}{\pi_\theta(x)} \right) \pi_*(x).$$

Based on training data, a discrete score matching estimator is a minimizer of

$$\sum_{k=1}^{N} \sum_{y \sim x_k} \left| 1 - \frac{\pi_\theta(y)}{\pi_\theta(x_k)} \right|^2 + 2 \sum_{k=1}^{N} \sum_{y \sim x_k} \left( \frac{\pi_\theta(x_k)}{\pi_\theta(y)} - \frac{\pi_\theta(y)}{\pi_\theta(x_k)} \right). \qquad (17.15)$$

$\blacklozenge$

## 17.3   Incomplete observations for graphical models

### 17.3.1   The EM Algorithm

Missing variable sin the context of graphical models may correspond to real processes that cannot be measured, which is common, for example, with biological data. They may be more conceptual objects that are interpretable but are not parts of the data acquisition process, like phonemes in speech recognition, or edges and labels in image processing and object recognition. They may also be variables that have been added to the model to increase its parametric dimension without increasing the complexity of the graph. However, as we will see, dealing with incomplete or imperfect observations brings the parameter estimation problem to a new level of difficulty.

Since it is the most common approach to address incomplete or noisy observations, we start with a description of how the EM algorithm (Algorithm 16.1) applies to graphical models, and of its limitations. We assume a graphical model on an undirected graph $G = (V, E)$, in which we assume that $V$ is separated in two non-intersecting subsets, $V = S \cup H$. Letting $X$ be a $G$-Markov random field, the part $X^{(S)}$ is assumed to be observable, and $X^{(H)}$ is hidden.

We assume that $X$ takes values in $\mathcal{F}(V)$, where we still denote by $F_s$ the sets in which $X_s$ takes values for $s \in V$. We let the model distribution belong to an exponential family, with

$$\pi_\theta(x) = \frac{1}{Z(\theta)} \exp\left(-\theta^T U(x)\right), \ x \in \mathcal{F}(V). \tag{17.16}$$

Assume that an $N$-sample $x_1^{(S)}, \ldots, x_N^{(S)}$ is observed over $S$. Since

$$\log \pi_\theta(x) = -\log Z(\theta) - \theta^T U(x),$$

the transition from $\theta_n$ to $\theta_{n+1}$ in Algorithm 16.1 is done by maximizing

$$-\log Z(\theta) - \theta^T \bar{U}_n \tag{17.17}$$

where

$$\bar{U}_n = \frac{1}{N} \sum_{k=1}^{N} E_{\theta_n}(U(X) \mid X^{(S)} = x_k^{(S)}). \tag{17.18}$$

So, the M-step of the EM, which maximizes (17.17), coincides with the complete-data maximum-likelihood problem for which the empirical average of $U$ is replaced by the average of its conditional expectations given the observations, as given in (17.18), which constitutes the E-step. As a consequence, a strict application of the EM algorithm for graphical models is unfeasible, since each step requires running an algorithm of similar complexity maximum likelihood for complete data, that we already identified as a challenging, computationally costly problem. The same remark holds for the SAEM algorithm of section 16.4.3, which also requires solving a maximum likelihood problem at each iteration.

### 17.3.2  Stochastic gradient ascent

The stochastic gradient ascent described in section 17.2.2 can be extended to partial observations [207], even though it loses the global convergence guarantee that resulted from the concavity of the log-likelihood for complete observations. Indeed, applying the computation of section 16.5.2, to a model given by (17.16), we get using proposition 17.4,

$$\partial_\theta \log \psi_\theta = E_\theta(E_\theta(U) - U \mid X^{(S)} = x^{(S)}) = E_\theta(U) - E_\theta(U \mid X^{(S)} = x^{(S)})$$

where we $\psi_\theta(x^{(S)})$ denotes the marginal distribution of $\pi_\theta$ on $S$.

Let $\pi_\theta(x^{(H)} \mid x^{(S)})$ denotes the conditional probability $P(X^{(H)} = x^{(H)} \mid X^{(S)} = s^{(S)})$ for the distribution $\pi_\theta$, therefore taking the form

$$\pi_\theta(x^{(H)} \mid x^{(S)}) = \frac{1}{\tilde{Z}(\theta, x^{(S)})} \exp\left(-\theta^T U(x^{(S)} \wedge x^{(H)})\right).$$

Assume given an ergodic transition probability $p_\theta$ on $\mathcal{F}(V)$, and a family of ergodic transition probabilities $p_\theta^{x^{(S)}}$, $x^{(S)} \in \mathcal{F}(S)$, such that the invariant distribution of $p_\theta$ is $\pi_\theta$, and the one of $p_\theta^{x^{(S)}}$ is $\pi_\theta(\cdot \mid x^{(S)})$. Then the following SGA algorithm can be used to estimate $\theta$

---

**Algorithm 17.1**

Start the algorithm with an initial parameter $\theta(0)$ and initial configurations $x(0)$ and $x_k^{(H)}(0)$, $k = 1,\ldots,N$. Then, at step $n$,

(SGH1)  Sample from the distribution $p_{\theta(n)}(x(n),\cdot)$ to obtain new configurations $x(n+1) \in \mathcal{F}(V)$.

(SGH2)  For $k = 1,\ldots,N$, sample from the distribution $p_{\theta(n)}^{x_k^{(S)}}(x_k^{(H)}(n),\cdot)$ to obtain a new configuration $x_k^{()}H(n+1)$ over the hidden vertexes.

(SGH3)  Update the parameter using

$$\theta(n+1) = \theta(n) + \gamma(n+1)\left( U(x(n+1)) - \frac{1}{N}\sum_{k=1}^{N} U(x_k^{(S)} \wedge x_k^{(H)}(n+1)) \right). \qquad (17.19)$$

---

### 17.3.3   Pseudo-EM Algorithm

The EM update

$$\theta_{n+1} = \underset{\theta}{\operatorname{argmax}}\left( \sum_{k=1}^{N} E_{\theta_n}\left( \log \pi_\theta(X) \mid X^{(S)} = x_k^{(S)} \right) \right).$$

being challenging for Markov random fields, it is tempting to replace the log-likelihood in the expectation by an other contrast, such as the log-pseudo-likelihood. A similar approach to that described here was introduced in Chalmond [51], for situations when the conditional distribution of $X^{(S)}$ given $X^{(H)}$ is "simple enough" (for example, if the variables $X_s, s \in S$ are conditionally independent given $X^{(H)}$) and when the cardinality of the sets $F_s$, $s \in H$ is small (binary, or ternary, variables).

The algorithm has the following variational interpretation. Fix $x^{(S)} \in \mathcal{F}(S)$ and $s \in H$. Also denote $\mu_s = 1/|\mathcal{F}(H \setminus \{s\})|$. If $q$ is a transition probability from $\mathcal{F}(H \setminus \{s\})$ to $F_s$, let

$$\Delta_\theta^{(s)}(q,x^{(S)}) = \sum_{y \in \mathcal{F}(H)} \left( \log\left( \frac{\pi_{\theta,s}(y^{(s)} \wedge x^{(S)} \mid y^{(H \setminus \{s\})})}{q(y^{(H \setminus \{s\})}, y^{(s)})\mu_s} \right) q(y^{(H \setminus \{s\})}, y^{(s)})\mu_s \right). \qquad (17.20)$$

This function is concave in $q$, since its first partial derivative with respect to $q(y^{(H\setminus\{s\})}, y^{(s)})$ (for each $y \in \mathcal{F}(H)$) is given by

$$\mu_s \log \pi_{\theta,s}(y^{(s)} \wedge x^{(S)} \mid y^{(H\setminus\{s\})}) \mu_s(y^{(H\setminus\{s\})}) - \mu_s \log(q(y^{(H\setminus\{s\})}, y^{(s)}) \mu_s) - \mu_s$$

so that its Hessian is the diagonal matrix with negative entries $-\mu_s/q(y^{(H\setminus\{s\})}, y^{(s)})$. Using Lagrange multipliers to express the constraints $\sum_{y^{(s)} \in F_s} q(y^{(H\setminus\{s\})}, y^{(s)}) = 1$ for all $y^{(H\setminus\{s\})}$, we find that $\Delta_\theta^{(s)}(q, x^{(S)})$ is maximized when $q(y^{(H\setminus\{s\})}, y^{(s)})$ is proportional to $\pi_{\theta,s}(y^{(s)} \wedge x^{(S)} \mid y^{(H\setminus\{s\})})$, yielding

$$q(y^{(H\setminus\{s\})}, y^{(s)}) = \pi_{\theta,s}(y^{(s)} \mid x^{(S)} \wedge y^{(H\setminus\{s\})}).$$

Now, consider the problem of maximizing

$$\sum_{k=1}^{N} \sum_{s \in H} \Delta_\theta^{(n)}(q_k^{(s)}, x_k^{(s)}) \tag{17.21}$$

with respect to $\theta$ and $q_k^{(s)}$, $k = 1, \ldots, N$, $s \in H$. Consider an iterative maximization scheme in which, from a current parameter $\theta_n$, one first, maximizes (17.21) with respect to transition probabilities $q_k^{(s)}$, then with respect to $\theta$ to obtain $\theta_{n+1}$. This scheme provides the iteration

$$\theta_{n+1} =$$

$$\operatorname*{argmax}_{\theta} \sum_{k=1}^{N} \sum_{s \in H} \sum_{y \in \mathcal{F}(H)} \left( \log \pi_{\theta,s}(y^{(s)} \wedge x_k^{(S)} \mid y^{(H\setminus\{s\})}) \right) \pi_{\theta_n,s}(y^{(s)} \mid x_k^{(S)} \wedge y^{(H\setminus\{s\})}) \mu_s.$$

### 17.3.4 Partially-observed Bayesian networks on trees

We now consider the situation in which the joint distribution of $X = X^{(S)} \wedge X^{(H)}$ is a Bayesian network over a directed acyclic graph $G = (V, E)$.

Assume that $x_1^{(S)}, \ldots, x_N^{(S)}$ are observed. The parameter $\theta$ is the collection of all $p(x^{(pa(s))}, x^{(s)})$ for $s \in V$. Define the random variables $I_{s,x}(y)$ equal to one if $y^{(\{s\} \cup pa(s))} = x^{(\{s\} \cup pa(s))}$ and zero otherwise. We can write

$$\log \pi(y) = \sum_{s \in S} \log p_s(y^{(pa(s))}, y^{(s)}) = \sum_{s \in S} \sum_{x^{(\{s\} \cup pa(s))} \in \mathcal{F}(\{s\} \cup pa(s))} \log p_s(x^{(pa(s))}, x^{(s)}) I_{s,x}(y)$$

This implies that

$$\sum_{k=1}^{N} E_{\theta_n}\left(\log \pi(x_k^{(S)}, X^{(H)}) \mid X^{(S)} = x_k^{(S)}\right)$$

$$= \sum_{x^{(\{s\}\cup pa(s))}\in\mathcal{F}(\{s\}\cup pa(s))} \log p_s(x^{(pa(s))}, x^{(s)}) \sum_{k=1}^{N} E_{\theta_n}(I_{s,x}(X) \mid X^{(S)} = x_k^{(S)})$$

$$= \sum_{x^{(\{s\}\cup pa(s))}\in\mathcal{F}(\{s\}\cup pa(s))} \log p_s(x^{(pa(s))}, x^{(s)}) \sum_{k=1}^{N} \pi_{\theta_n}(x^{(\{s\}\cup pa(s))} \mid X^{(S)} = x_k^{(S)}).$$

The EM iteration at step $n$ then is

$$p_s^{(n+1)}(x^{(pa(s))}, x^{(s)}) = \frac{1}{Z_s(x^{(s^-)})} \sum_{k=1}^{N} \pi_{\theta_n}(x^{(\{s\}\cup pa(s))} \mid X^{(S)} = x_k^{(S)})$$

with

$$\pi_{\theta_n}(x) = \prod_{s\in V} p^{(n)}(x^{(pa(s))}, x^{(s)}),$$

$Z_s$ being a normalization constant.

If the estimation is solved with a Dirichlet prior $\text{Dir}(1 + a_s(x^{(s)}, x^{(pa(s))}))$, the update formula becomes

$$p_s^{(n+1)}(x^{(pa(s))}, x^{(s)}) = \frac{1}{Z_s(x^{(s^-)})}\left(a_s(x^{(s)}, x^{(pa(s))}) + \sum_{k=1}^{N} \pi_{\theta_n}(x^{(\{s\}\cup pa(s))} \mid X^{(S)} = x_k^{(S)})\right).$$

$$(17.22)$$

This algorithm is very simple when the conditional distributions $\pi_{\theta_n}(x^{(s\cup pa(s))} \mid X^{(S)} = x_k^{(S)})$ can be easily computed, which is not always the case for a general Bayesian network, since conditional distributions do not always have a structure of Bayesian network. The computation is simple enough for trees, however, since conditional tree distributions are still trees (or forests). More precisely, the conditional distribution given the observed variables can be written in the form

$$\pi(y^{(H)} \mid x^{(S)}) = \frac{1}{Z(x^{(S)})} \prod_{s\in H} \varphi_{s,x}(y^{(s)}) \prod_{t\sim s,\{s,t\}\subset H} \varphi_{st}(y^{(s)}, y^{(t)})$$

with $\varphi_{s,pa(s)}(y^{(s)}, y^{(pa(s))}) = p_s(y^{(pa(s))}, y^{(s)})$ and, letting $\varphi_s(y^{(s)}) = p_s(y^{(s)})$ if $pa(s) = \emptyset$ and 1 otherwise,

$$\varphi_{s,x}(y^{(s)}) = \varphi_s(y^{(s)}) \prod_{t\sim s,t\in S} \varphi_{st}(y^{(s)}, x^{(t)}).$$

So, the marginal joint distribution of a vertex and its parents are directly given by belief propagation, using the just defined interactions. This training algorithm is summarized below.

---

**Algorithm 17.2 (Learning tree distributions with hidden variables)**
Start with some initial guess of the conditional probabilities (for example, those given by the prior). The iterate the following two steps providing the transition from $\theta_n$ to step $\theta_{n+1}$.

(1) For $k = 1, \ldots, N$, use belief propagation (or sum-prod) to compute all $\pi_{\theta_n}(x^{(\{s\} \cup pa(s))} \mid X^{(S)} = x_k^{(S)})$. Note that these probabilities can be 0 or 1 when $s \in S$ and/or $pa(s) \subset S$.

(2) Use (17.22) to compute the next set of parameters.

---

The tree case includes the important example of hidden Markov models, which are defined as follows. $S$ and $H$ are ordered, with same cardinality, say $S = \{s_1, \ldots, s_q\}$ and $H = \{h_1, \ldots, h_q\}$. Edges are $(h_1, h_2), \ldots, (h_{q-1}, h_q)$ and $(h_1, s_1), \ldots, (h_q, s_q)$. The interpretation generally is that the hidden variables, $h_s$, are the variables of interest, and behave like a Markov chain, and that the observations, $x_s$, are either noisy or transformed versions of them. A major application is in speech recognition, where the $h_s$'s are labels that represent specific phonemes (little pieces of spoken words) and the $x_s$'s are measured signals. The transitions between hidden variables then describe how phonemes are likely to appear in sequence for a given language, and those between hidden and observed variables describe how each phoneme is likely to be pronounced and heard.

### 17.3.5   General Bayesian networks

The algorithm in the general case can move from tractable to intractable depending on the situation. This must generally be handled in a case by case basis, by analyzing the conditional structure, for a given model, knowing the observations.

In practice, it is always possible to use loopy belief propagation to obtain some approximation of the conditional probabilities, even if it is not sure that the algorithm will converge to the correct marginals. When feasible, junction trees can be used, too. Monte-Carlo sampling is also an option, although quite computational.

# Chapter 18

# Deep Generative Methods

## 18.1 Normalizing flows

### 18.1.1 General concepts

We develop, in this chapter, methods that model stochastic processes using a feed-forward approach that generates complex random variables using non-linear transformations of simpler ones. Many of these methods can be seen as instances of structural equation models (SEMs), described in section 15.3, with, for deep-learning implementations, high-dimensional parametrizations of (15.8).

With start with the formally simple case where the modeled variable takes values in $\mathbb{R}^d$ and is modeled as

$$X = g(Z)$$

where $Z$ also takes values in $\mathbb{R}^d$, with a known distribution and $g$ is $C^1$, invertible, with a $C^1$ inverse on $\mathbb{R}^d$, i.e., is a *diffeomorphism* of $\mathbb{R}^d$. Let us denote by $h$ the inverse of $g$.

If $Z$ has a p.d.f. $f_Z$ with respect to Lebesgue's measure, then, using the change of variable formula, the p.d.f. of $X$ is

$$f_X(x) = f_Z(h(x)) \, |\det \partial_x h(x)|.$$

Now, given a training set $T = (x_1, \ldots, x_N)$, the log-likelihood, considered as a function of $h$, is given by

$$\ell(h) = \sum_{k=1}^{N} \log f_Z(h(x_k)) + \sum_{k=1}^{N} \log |\det \partial_x h(x_k)|. \tag{18.1}$$

This expression should then be maximized with respect to $h$, subject to some restrictions or constraints to avoid overfitting.

### 18.1.2   A greedy computation

One can define a rich class of diffeomorphisms through iterative compositions of simple transformations. This framework was introduced in [187], where a greedy approach was suggested to build such compositions. The method was termed "normalizing flows," since it create a discrete flow of diffeomorphisms that transform the data into a sample of a normal distribution.

We quickly describe the basic principles of the algorithm. One starts with a parametrized family, say $(\psi_\alpha, \alpha \in A)$ of diffeomorphisms of $\mathbb{R}$. Such families are relatively easy to design, one example proposed in [187] being a smoothed version of the piecewise linear function

$$u \mapsto v_0 + (1 - \sigma)u + \gamma|(1 - \sigma)u - u_0|$$

which is increasing as soon as $0 \leq \max(\sigma, \gamma) < 1$. The smoothed version has an additional parameter, $\epsilon$, and takes the form

$$u \mapsto v_0 + (1 - \sigma)u + \gamma\sqrt{\epsilon^2 + ((1 - \sigma)u - u_0)^2}.$$

This transformation is parametrized by $\alpha = (v_0, \sigma, \gamma, u_0, \epsilon)$. Other families of parametrized transformations can be designed. A multivariate transformation $\varphi_{\boldsymbol{\alpha}, U} : \mathbb{R}^d \to \mathbb{R}^d$ can then be associated to families $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_d)$ and orthogonal matrices $U$ by taking

$$\varphi_{\boldsymbol{\alpha}, U}(x) = \begin{pmatrix} \psi_{\alpha_1}(y^{(1)}) \\ \vdots \\ \psi_{\alpha_d}(y^{(d)}) \end{pmatrix}$$

with $y = Ux$.

The algorithm in [187] is initialized with $h_0 = \mathrm{id}\,[d]$ and update the transformation at step $n$ according to

$$h_n = \varphi_{\boldsymbol{\alpha}_n, U_n} \circ h_{n-1}.$$

In this update, $U_n$ is generated as a random rotation matrix, and $\boldsymbol{\alpha}_n$ is determined as a gradient ascent update (starting from $\boldsymbol{\alpha} = 0$) for the maximization of

$$\boldsymbol{\alpha} \mapsto \ell(\varphi_{\boldsymbol{\alpha}, U_n} \circ h_{n-1}).$$

(Here, the current value $h_{n-1}$ is not revisited, therefore providing a "greedy" optimization method.)

Letting $z_{n,k} = h_n(x_k)$, the chain rule implies that

$$\ell(\varphi_{\alpha,U_n} \circ h_{n-1})) = \sum_{k=1}^{N} \log f_Z(\varphi_{\alpha,U_n}(z_{n-1,k})) + \sum_{k=1}^{N} \log|\det \varphi_{\alpha,U_n}(z_{n-1,k})|$$
$$+ \sum_{k=1}^{N} \log|\det \partial_x h_{n-1}(x_k)|.$$

Since the last term does not depend on $\alpha$, we see that it suffices to keep track of the "particle" locations, $z_{n-1,k}$ to be able to compute $\alpha_n$. Note also that these locations are easily updated with $z_{n,k} = \varphi_{\alpha_n,U_n}(z_{n-1,k})$.

### 18.1.3   Neural implementation

This iterated composition of diffeomorphisms obviously provides a neural architecture similar to those discussed in chapter 11. Fixing the number of iterations to be, say, $m$, one can consider families of diffeomorphisms $(\varphi_\theta)$ indexed by a parameter $w$ (we had $w = (\alpha, U)$ in the previous discussion), and optimize (18.1) over all functions $h$ taking the form $h = \varphi_{w_m} \circ \cdots \circ \varphi_{w_1}$. Letting $z_{j,k} = \varphi_{w_j} \circ \cdots \circ \varphi_{w_1}(x_k)$ for $j \leq m$ (with $z_{0,k} = x_k$), we can write

$$\ell(h) = \sum_{k=1}^{N} \log f_Z(z_{m,k}) + \sum_{k=1}^{N}\sum_{j=1}^{m} \log|\det \partial_x \varphi_{w_j}(z_{j-1,k})|.$$

Normalizing flows in this form are described in [161, 107, 148]. The gradient of $\ell$ with respect to the parameters $w_1,\ldots,w_m$ can be computed by backpropagation. We note however that, unlike typical neural implementations, the parameters may come with specific constraints, such as $U \in \mathcal{O}_d(\mathbb{R})$ when $w = (\alpha, U)$, so that the gradient and associated displacement may have to be adapted compared to standard gradient ascent implementations (see section 20.6.3 for a discussion of first-order implementations of gradient methods for functions of orthogonal matrices, and [1] for more general methods on optimization over matrix groups).

### 18.1.4   Time-continuous version

In section 11.6, we described how diffeomorphisms could be generated as flows of differential equations, and this remark can be used to provide a time-continuous version of normalizing flows. Using (11.3), one generates trajectories $z(\cdot)$ by solving over, say, $[0, T]$

$$\partial_t z(t) = \psi_{w(t)}(z(t))$$

with $z(0) = x$ for some function $w : t \mapsto w(t)$. Letting $z(t) = h_w(t,x)$ (which defines $h_w$), we know that, under suitable assumptions on $\psi$, the mapping $x \mapsto h_w(t,x)$ is a diffeomorphism of $\mathbb{R}^d$. One can then maximize

$$\ell(h_w(T, \cdot)) = \sum_{k=1}^{N} \log f_Z(h_w(T, x_k)) + \sum_{k=1}^{N} \log|\det \partial_x h_w(T, x_k)|$$

with respect to the function $w$. Let $z_k(t) = h_w(t, x_k)$ and $J_k(t) = \log|\det \partial_x h_w(t, x_k)|$. We have, by definition

$$\partial_t z_k(t) = \psi_{w(t)}(z_k(t))$$

with $z_k(0) = x_k$. One can also show that

$$\partial_t J_k(t) = \nabla \cdot \psi_{w(t)}(z_k(t))$$

with $J_k(0) = 0$, where the r.h.s. is the divergence of $\psi_{w(t)}$ evaluated at $z_k(t)$. We provide a quick (and formal) justification of this fact. First note that differentiating $\partial_t h_w(t, x) = \psi_{w(t)}(h_w(t, x))$ with respect to $x$ yields

$$\partial_t \partial_x h_w(t, x) = \partial_x \psi_{w(t)}(h_w(t, x)) \partial_x h_w(t, x).$$

The mapping $\mathcal{J} : A \mapsto \log|\det(A)|$ is differentiable on the set of invertible matrices and is such that $d\mathcal{J}(A)H = \text{trace}(A^{-1}H)$. Applying the chain rule, we find

$$\partial_t \log|\det \partial_x h_w(t, x)| = \text{trace}(\partial_x h_w(t, x)^{-1} \partial_x \psi_{w(t)}(h_w(t, x)) \partial_x h_w(t, x))$$
$$= \text{trace}(\partial_x \psi_{w(t)}(h_w(t, x))) = \nabla \cdot \psi_{w(t)}(h_w(t, x)).$$

From this, it follows that the time-continuous normalizing flow problem can be reformulated as maximizing

$$\sum_{k=1}^{N} \log f_Z(z_k(T)) + \sum_{k=1}^{N} J_k(T)$$

subject to $\partial_t z_k(t) = \psi_{w(t)}(z_k(t))$, $\partial_t J_k(t) = \nabla \cdot \psi_{w(t)}(z_k(t))$, $z_k(0) = x_k$, $J_k(0) = 0$. This is an optimal control problem, whose analysis can be done similarly to that made in section 11.6.1, provided that $\nabla \cdot \psi_{w(t)}$ can be expressed in closed form.

Note that the inverse of $h_w(T, \cdot)$, which provides the generative model going from $Z$ to $X$ can also be obtained as the solution of an ODE. Namely, if one solves the differential equation

$$\partial_t x(t) = -\psi_{w(T-t)}(x(t))$$

with initial condition $x(0) = z$, then $x(T)$ solves the equation $h_w(T, \cdot) = z$.

## 18.2 Non-diffeomorphic models and variational autoencoders

### 18.2.1 General framework

The previous discussion addressed the situation $X = g(Z)$ when $g$ is a diffeomorphism, which required, in particular, that $X$ and $Z$ are real vectors with identical dimensions. This may not always be desirable, as one may prefer a small-dimensional variable $Z$ (in the spirit of the factor analysis methods discussed in chapter 20), or a high-dimensional $Z$ to increase, for example the modeling power. In addition, the observation variables may be discrete, which precludes the use of the change of variables formula. In such cases, $Z$ has to be treated as a hidden variable using one of the methods discussed in chapter 16.

It will convenient to model the generative process in the form of a conditional distribution of $X$ given $Z$ rather than a deterministic function. We place ourselves in the framework of chapter 16 (with slightly modified notation) and let $\mathcal{R}_X$ and $\mathcal{R}_Z$ denote the measured spaces over where $X$ and $Z$ take their values, with measures $\mu_X$ and $\mu_Z$, and assume that the conditional distribution of $X$ given $Z = z$ has density $f_X(x \mid z, \theta)$ with respect to $\mu_X$, for some parameter $\theta$. We also assume that $Z$ has a distribution with density $f_Z$ with respect to $\mu_Z$, that we assume given and unparametrized. One can then directly apply the algorithms provided in chapter 16, and in particular the variational methods described in section 16.4.4 with an appropriate definition of the approximation of the conditional density of $Z$ given $X$. An important example in this context is provided by variational autoencoders (VAEs) that we now present.

### 18.2.2 Generative model for VAEs

VAEs [103, 104] model $X \in \mathbb{R}^d$ as $X = g(Z, \theta) + \epsilon$ where $\epsilon$ is a centered Gaussian noise with covariance matrix $Q$. The function $g$ is typically non-linear, and VAEs have been introduced with this function modeled as a deep neural network (see chapter 11). Letting $\varphi_{\mathcal{N}}(\cdot; 0, Q)$ denote the p.d.f. of the Gaussian distribution $\mathcal{N}(0, Q)$, the conditional distribution of $X$ given $Z = z$ has density

$$f_X(x \mid z, \theta) = \varphi_{\mathcal{N}}(x - g(z, \theta)); 0, Q)$$

with respect to Lebesgue's measure on $\mathbb{R}^d$.

Following the procedure in section 16.4.4, we define an approximation of the conditional distribution of $Z$ given $X$. Assuming that $Z \in \mathbb{R}^p$, we let this distribution be $\mathcal{N}(\mu(x, w), \Sigma(x, w))$ for some functions $\mu$ and $\Sigma$, $w$ being a parameter. To ensure that $\Sigma \geq 0$, we will represent it in the form $\Sigma(x, w) = S(x, w)^2$ where $S$ is a symmetric matrix. In [103], both functions $\mu$ and $S$ are represented as neural networks

parametrized by $w$. The joint density of $X$ and $Z$ is such that

$$\log f_{X,Z}(x,z;\theta,Q) = \log \varphi_{\mathcal{N}}(x - g(z,\theta)); 0, Q) + \log f_Z(z)$$

$$= -\frac{1}{2}(x - g(z,\theta))^T Q^{-1}(x - g(z,\theta)) - \frac{1}{2}\log \det Q - \frac{d}{2}\log 2\pi + \log f_Z(z)$$

We also have

$$\log \varphi_{\mathcal{N}}(z; \mu(x,w), S(x,w)^2) = -\frac{1}{2}(z - \mu(x,w))^T S(x,w)^{-2}(z - \mu(x,w)) - \log \det S(x,w) - \frac{p}{2}\log 2\pi.$$

We can then rewrite the algorithm in (16.15) as

$$\begin{cases} \theta_{n+1} = \theta_n + \gamma_{n+1}\partial_\theta \log f_{X,Z}(X_{n+1}, Z_{n+1}; \theta_n, Q_n) \\ Q_{n+1} = Q_n + \gamma_{n+1}\partial_Q \log f_{X,Z}(X_{n+1}, Z_{n+1}; \theta_n, Q_n) \\ w_{n+1} = w_n + \gamma_{n+1}\log\left(\dfrac{f_{X,Z}(X_{n+1}, Z_{n+1}; \theta_n, Q_n)}{\varphi_{\mathcal{N}}(X_{n+1}; \mu(X_{n+1}, w_n), S(X_{n+1}, w_n)^2)}\right) \\ \qquad\qquad \times \partial_w \log \varphi_{\mathcal{N}}(X_{n+1}; \mu(X_{n+1}, w_n), S(X_{n+1}, w_n)^2) \end{cases} \quad (18.2)$$

where $X_{n+1}$ is drawn uniformly from the training data and

$$Z_{n+1} \sim \mathcal{N}(\mu(X_{n+1}, w_n), S(X_{n+1}, w_n)^2).$$

The derivatives in this system can be computed from those of $g, \mu$ and $S$ (typically involving back-propagation) and the expression of the derivatives of the determinant and inverse of a matrix provided in (1.4) and (1.6).

The computation can be simplified if one assumes that $f_Z$ is the p.d.f. of a standard Gaussian, i.e., $f_Z = \varphi_{\mathcal{N}}(\cdot; 0, \mathrm{Id}_{\mathbb{R}^p})$. Indeed, in that case, the integral in (16.11), which is, using the current notation

$$\int_{\mathbb{R}^p} \log \frac{\varphi_{\mathcal{N}}(x - g(z,\theta); 0, Q)\varphi_{\mathcal{N}}(z; 0, \mathrm{Id}_{\mathbb{R}^p})}{\varphi_{\mathcal{N}}(z; \mu(x,w), S(x,w)^2)} \varphi_{\mathcal{N}}(z; \mu(x,w), S(x,w)^2) dz, \quad (18.3)$$

can be partially computed. For any two $p$-dimensional Gaussian p.d.f.'s, one has

$$\int_{\mathbb{R}^p} \log \varphi_{\mathcal{N}}(z; \mu_1, \Sigma_1) \, \varphi_{\mathcal{N}}(z; \mu_2, \Sigma_2) \, dz = -\frac{1}{2}\mathrm{trace}(\Sigma_1^{-1}\Sigma_2) - \frac{1}{2}(\mu_2 - \mu_1)^T \Sigma_1^{-1}(\mu_2 - \mu_1)$$

$$-\frac{1}{2}\log \det(\Sigma_1) - \frac{p}{2}\log(2\pi). \quad (18.4)$$

As a consequence, (18.3) becomes

$$-\frac{1}{2}\mathbb{E}_w\left((X - g(Z,\theta))^T Q^{-1}(X - g(Z,\theta))\right) - \frac{1}{2}\log \det Q - \frac{d}{2}\log 2\pi$$

$$-\mathbb{E}_w\left(\frac{1}{2}\mathrm{trace}(S(X,w)^2) + \frac{1}{2}|\mu(X,w)|^2 - \log \det(S(X,w))\right) + \frac{p}{2}, \quad (18.5)$$

where $\mathbb{E}_w$ denotes the expectation for the random variable $(X, Z)$ where $X$ follows a uniform distribution over training data and the conditional distribution of $Z$ given $X = x$ is $\mathcal{N}(\mu(x, w), S(x, w)^2)$.

The algorithm proposed in Kingma and Welling [103] introduces a change of variable $Z = \mu(X, w) + S(X, w)U$ where $U \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^p})$, rewriting (18.5) as

$$-\frac{1}{2}\mathbb{E}\left((X - g(\mu(X, w) + S(X, w)U, \theta))^T Q^{-1}(X - g(\mu(X, w) + S(X, w)U, \theta))\right)$$

$$-\mathbb{E}_w\left(\frac{1}{2}\mathrm{trace}(S(x, w)^2) + \frac{1}{2}|\mu(X, w)|^2 - \log\det(S(X, w))\right)$$

$$-\frac{1}{2}\log\det Q - \frac{d}{2}\log 2\pi + \frac{p}{2}, \quad (18.6)$$

with a modified version of (18.2). Letting

$$F(\theta, Q, w, x, u) = -\frac{1}{2}(x - g(\mu(x, w) - S(x, w)U, \theta))^T Q^{-1}(x - g(\mu(x, w) - S(x, w)U, \theta))$$

$$-\frac{1}{2}\log\det Q - \frac{1}{2}\mathrm{trace}(S(x, w)^2) - \frac{1}{2}|\mu(x, w)|^2 + \log\det(S(x, w))$$

the resulting algorithm is

$$\begin{cases} \theta_{n+1} = \theta_n + \gamma_{n+1}\partial_\theta F(\theta_n, Q_n, w_n, X_{n+1}, U_{n+1}) \\ Q_{n+1} = Q_n + \gamma_{n+1}\partial_Q F(\theta_n, Q_n, w_n, X_{n+1}, U_{n+1}) \\ w_{n+1} = w_n - \gamma_{n+1}\partial_w F(\theta_n, Q_n, w_n, X_{n+1}, U_{n+1}) \end{cases} \quad (18.7)$$

where $X_{n+1}$ is drawn uniformly from the training data and $U_{n+1} \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^p})$.

### 18.2.3  Discrete data

This framework can be easily adapted to situations in which the observations are discrete. Consider, as an example, the situation in which $X$ takes values in $\{0, 1\}^V$, where $V$ is a set of vertexes, i.e., $X$ is a binary Markov random field on $V$. Assume, as a generative model, that conditionally to the latent variable $Z \in \mathbb{R}^p$, the variables $X^{(s)}, s \in V$ are independent and $X^{(s)}$ follows a Bernoulli distribution with parameter $g_s(z, \theta)$, where $g : \mathbb{R}^p \to [0, 1]^V$. Assume also that $Z \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^p})$, and define, as above, an approximation of the conditional distribution of $Z$ given $X = x$ as a Gaussian with mean $\mu(x, w)$ and covariance matrix $S(x, w)^2$. Then, the joint density of $X$ and $Z$ (with respect to the product of the counting measure on $\{0, 1\}^V$ and Lebesgue's measure on $\mathbb{R}^p$) is

$$\log f_{X,Z}(x, z; \theta) = x \log g(z, \theta) + (1 - x)\log(1 - g(z, \theta)) + \log \varphi_{\mathcal{N}}(z; 0, \mathrm{Id}_{\mathbb{R}^p})$$

and (18.2) becomes

$$
\begin{cases}
\theta_{n+1} = \theta_n + \gamma_{n+1} \partial_\theta \log f_{X,Z}(X_{n+1}, Z_{n+1}; \theta_n, Q_n) \\
\\
w_{n+1} = w_n + \gamma_{n+1} \log\left( \dfrac{f_{X,Z}(X_{n+1}, Z_{n+1}; \theta_n)}{\varphi_{\mathcal{N}}(X_{n+1}; \mu(X_{n+1}, w_n), S(X_{n+1}, w_n)^2)} \right) \\
\qquad\qquad \times \partial_w \log \varphi_{\mathcal{N}}(X_{n+1}; \mu(X_{n+1}, w_n), S(X_{n+1}, w_n)^2)
\end{cases}
\tag{18.8}
$$

## 18.3   Generative Adversarial Networks (GAN)

### 18.3.1   Basic principles

Similarly to the methods discussed so far, GANs [82], use a one-step nonlinear generator $X = g(Z, \theta)$, with $\theta \in \mathbb{R}^K$, to model observed data (we here switch back to a deterministic relation), where $Z$ has a known distribution, with p.d.f. $f_Z$, for example $Z \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^p})$. However, unlike the exact or approximate likelihood maximization that were discussed in sections 18.1 and 18.2, GANs us a different criterion for estimating the parameter $\theta$ by minimizing metrics that can be approximated by optimizing a classifier. The classifier is a function $x \mapsto f(x, w)$, parametrized by $w \in \mathbb{R}^M$, whose goal is to separate simulated samples from real ones: it takes values in $[0, 1]$ and estimates the (posterior) probability that its input $x$ is real. GANs' adversarial paradigm consists in estimating $\theta$ and $w$ together so that generated data, using $\theta$, are indistinguishable from real ones using the optimal $w$. Their basic structure is summarized in Figure 18.1.



Figure 18.1:  Basic structure of GANs: $W$ is optimized to improve the prediction problem: "real data" vs. "simulation". Given $W$, $\theta$ is optimized to worsen the prediction.

### 18.3.2   Objective function

Let $P_\theta$ denote the distribution of $g(Z, \theta)$, and $P_{\text{true}}$ the target distribution of "real data." One can formalize the "real data" vs. "simulation" problem with a pair of

random variables $(X, Y)$ where $Y$ follows a Bernoulli distribution with parameter $1/2$, and the conditional distribution of $X$ given $Y$ is $P_{\text{true}}$ when $Y = 1$ and $P_\theta$ when $Y = 0$. Given a loss function $r : \{0, 1\} \times [0, 1] \to [0, +\infty)$, one can define

$$U(\theta, w) = E_\theta(r(Y, f(X, w)))$$

and

$$U^*(\theta) = \min_{w \in \mathbb{R}^M} U(\theta, w).$$

We want to maximize $U^*$ or, equivalently, solve the optimization problem

$$\theta^* = \text{argmax}_\theta \min_{w \in \mathbb{R}^M} U(\theta, w).$$

Note that

$$2U(\theta, w) = E_{\text{true}}(r(1, f(X, w))) + E_\theta(r(0, f(X, w)))$$

so that choosing the cost requires to specify the two functions $t \mapsto r(1, t)$ and $t \mapsto r(0, t)$. In Goodfellow et al. [82], they are:

$$\begin{aligned} r(1, t) &= -\log t \\ r(0, t) &= -\log(1 - t). \end{aligned} \tag{18.9}$$

### 18.3.3 Algorithm

Using costs in (18.9), one must compute

$$\theta^* = \text{argmin} \max_{w \in \mathbb{R}^M} \left( E_{\text{true}}(\log f(X, w)) + E_\theta(\log(1 - f(X, w))) \right)$$

$$= \text{argmin} \max_{w \in \mathbb{R}^M} \left( E_{\text{true}}(\log f(X, w)) + E(\log(1 - f(g(Z, \theta), w))) \right).$$

Such min-max, or saddle-point problem are numerically challenging. The following algorithm was proposed in Goodfellow et al. [82], and also includes a stochastic approximation component. Indeed, in practice, $E_{\text{true}}$ is only known through the observation of training data, say $x_1, \dots, x_N$. Moreover, $E_\theta$ is only accessible through Monte-Carlo simulation, so that both expectations can only be approximated through finite-sample averaging.

---

**Algorithm 18.1 (GAN training algorithm)**

1. Extract a batch of $m$ examples from training data, simulate $m$ samples according to $P_\theta$ and run a few (stochastic) gradient ascent steps with fixed $\theta$ to update $w$, replacing expectations by averages.

2. Generate $m$ new samples of $Z$ and update $\theta$ with fixed $w$ by iterating a few steps of (stochastic) gradient descent.

---

### 18.3.4   Associated probability metric and Wasserstein GANs

Let $\mathcal{F}$ be the family of all measurable functions: $f : \mathbb{R}^d \to [0,1]$. Given two possible probability distributions $P_1, P_2$ (with associated expectations denoted $E_1, E_2$) of a random variable $X$ taking values in $\mathbb{R}^d$, consider the function

$$D(P_1, P_2) = 2\log 2 + \max_{f \in \mathcal{F}} \Big( E_1(\log f(X)) + E_2(\log(1 - f(X))) \Big)$$

Assume that $X$ under $P_1$ (resp. under $P_2$) has p.d.f. $g_1$ (resp. $g_2$) with respect to Lebesgue's measure (this assumption is not needed for the following to hold, but makes the discussion more elementary). Then

$$E_1(\log f(X)) + E_2(\log(1 - f(X))) = \int_{\mathbb{R}^d} (g_1 \log f + g_2 \log(1 - f)) dx$$

which is maximal at $f^* = g_1/(g_1 + g_2)$. For this $f^*$,

$$
\begin{aligned}
2\log 2 + E_1(\log f^*(X)) + E_2(\log(1 - f^*(X))) &= \int_{\mathbb{R}^d} g_1 \log \frac{2g_1}{g_1 + g_2} dx + \int_{\mathbb{R}^d} g_2 \log \frac{2g_2}{g_1 + g_2} dx \\
&= KL\left(\frac{g_1 + g_2}{2}, g_1\right) + KL\left(\frac{g_1 + g_2}{2}, g_2\right)
\end{aligned}
$$

This expression is called the *Jensen-Shannon divergence* between $g_1$ and $g_2$. It is always non-negative, and vanishes only when $g_1 = g_2$.

So, $D : (P_1, P_2) \mapsto D(P_1, P_2)$ can be interpreted as a way to evaluate the difference between two probability distributions on $\mathbb{R}^d$. One can then define

$$\hat{D}(P_1, P_2) = \max_{w \in \mathbb{R}^M} \Big( E_1(\log f(X, w)) + E_2(\log(1 - f(X, w))) \Big)$$

as an approximation of $D$ in which the set of all possible functions with values in $[0,1]$ is replaced by those arising from the GAN classification network, parametrized by $w$. This approximation is useful when $g_1, g_2$ are only observable through random sampling or simulation. With this interpretation, GANs minimize $\hat{D}(P_{\text{true}}, P_\theta)$.

This discussion suggests that new types of GAN may be designed using other discrepancy functions between probability distributions, provided they can be expressed in terms of the maximization of some quantity over some space of functions. Consider, for example the norm in total variation, defined by (for discrete distributions)

$$D_{\text{var}}(P_1, P_2) = \frac{1}{2} \sum_x |P_1(x) - P_2(x)|.$$

or, in the general case $D_{\text{var}}(P_1, P_2) = \max_A (P_1(A) - P_2(A))$.

If $\mathcal{F}$ is the space of continuous functions $f : \mathbb{R}^d \to [0,1]$, then we also have (under mild assumptions on $P_1$ and $P_2$)

$$D_{\text{var}}(P_1, P_2) = \max_{f \in \mathcal{F}}(E_1(f) - E_2(f)).$$

Since neural nets typically generate continuous functions with values in $[0,1]$, one could train GANs by maximizing

$$\hat{D}_{\text{var}}(P_1, P_2) = \max_{w \in \mathbb{R}^M}\Big(E_1(f(X,w)) - E_2(f(X,w))\Big)$$

However, the total variation distance is too crude to allow for meaningful comparisons between distributions. For example, the distance between two Dirac distributions at, say, $x_1$ and $x_2$ in $\mathbb{R}^d$ is always 1, whatever the distance between $x_1$ and $x_2$, unless $x_1 = x_2$. A more sensitive distance can be defined based on the notion of optimal transport.

The Monge-Kantorovich, also called Wasserstein, and sometimes also called "earth-mover", distance evaluates the minimal total distance along which "mass" needs to be transported to transform a distribution, $P_1$, into another, $P_2$. Its mathematical definition is

$$D_w(P_1, P_2) = \inf_Q \int_{\mathbb{R}^d \times \mathbb{R}^d} |x_1 - x_2| Q(dx_1, dx_2)$$

where the inf is computed over all joint distributions on $\mathbb{R}^d \times \mathbb{R}^d$ whose first marginal is $P_1$ and second marginal $P_2$. Note that the distance $D_w$ between $\delta_{x_1}$ and $\delta_{x_2}$ now is $|x_1 - x_2|$.

The Wasserstein distance can also be defined by

$$D_w(P_1, P_2) = \max_{f \in \mathcal{F}}(E_1(f) - E_2(f))$$

where $\mathcal{F}$ is now the space of contractive (or 1-Lipschitz) functions, i.e., $f \in \mathcal{F}$ if and only if, for all $x_1, x_2 \in \mathbb{R}^d$, $|f(x_1) - f(x_2)| \leq |x_1 - x_2|$.

Using the fact that a neural network with all weights bounded by a constant $K$ generates a function whose Lipschitz constant is controlled solely by $K$, one can then approximate (up to a multiplicative constant) the Wasserstein distance by

$$\hat{D}_w(P_1, P_2) = \max_{w \in \mathbb{W}}\Big(E_1(f(X,w)) - E_2(f(X,w))\Big)$$

where $\mathbb{W}$ is the set of all weights bounded by a fixed constant. Given the distribution $P_{\text{true}}$ and the model $P_\theta$, Wasserstein GANs (WGANs [11]) must then solve the saddle-point problem

$$U(\theta, w) = \max_{w \in \mathbb{W}}\Big(E_{\text{true}}(f(X,w)) - E_\theta(f(X,w))\Big)$$

and

$$U^*(\theta) = \min_{w \in \mathbb{R}^M} U(\theta, w),$$

with an algorithm similar to that described earlier.

As a final reference, we note the improved WGAN algorithm introduced in Gulrajani et al. [84] in which the boundedness constraint in the weights is replaced by an explicit control of the derivative in $x$ of the function $f$. More precisely, introduce a random variable $Z$ with distribution $\tilde{P}_\theta$ equal $(1-U)X + UX'$ where $U$ is uniformly distributed over $[0,1]$ and $X$ and $X'$ are independent respectively following the distribution $P_{\text{true}}$ and $P_\theta$. Then, the following approximation of the Wasserstein distance between $P_{\text{true}}$ and $P_\theta$ that is used in Gulrajani et al. [84]:

$$\hat{D}_w(P_{\text{true}}, P_\theta) = \max_{w \in \mathbb{W}} \left( E_{\text{true}}(f(X,w)) - E_\theta(f(X,w)) - \tilde{E}_\theta((|\partial_z f(Z,w)| - 1)^2) \right).$$

## 18.4   Reversed Markov chain models

### 18.4.1   General principles

The discussions in sections 18.2 and 18.3 can be applied to sequences of structural equations (describing finite Markov chains) in the form

$$\begin{cases} Z_0 = \xi_0 \\ Z_{k+1} = g(Z_k, \xi_k; \theta_k), \ k = 0, \dots, m-1 \\ \quad X = Z_m \end{cases}$$

where $\xi_0, \dots, \xi_{m-1}$ are random variables with fixed distribution.

Indeed, letting $\tilde{Z} = (\xi_0, \dots, \xi_{n-1})$ and $\tilde{\theta} = (\theta_0, \dots, \theta_{m-1})$ the whole system can be considered as a function $X = G(\tilde{Z}, \tilde{\theta})$ as considered in these sections. This representation, however, includes a large number of hidden variables, and it is unclear whether much improvement can be added to the case $m = 1$ to justify the additional computational load.

Reversed Markov chain models use a different generative approach in that they first model a forward Markov chain $Z_n, n \geq 0$ which is ergodic with known (and easy to sample from) limit distribution $Q_\infty$, and initial distribution $Q_{\text{true}}$, the true distribution of the data. If one fixes a large enough number of steps, say, $\tau$, then it is reasonable to assume that $Z_\tau$ approximately follows the limit distribution, $Q_\infty$. One can then (approximately) sample from $Q_{\text{true}}$ by sampling $\tilde{Z}_0$ according to $Q_\infty$ and then applying $\tau$ steps of the time-reversed Markov chain.

Reversed chains were discussed in section 12.3.3. Assuming that $Q_{\text{true}}$ and $P(z, \cdot)$ have a density with respect to a fixed measure $\mu$ on $\mathcal{R}_Z$, we found that $\tilde{Z}_k = Z_{\tau - k}$ is a non-homogeneous Markov chain whose transition probability $\tilde{P}_k(x, A) = P(\tilde{Z}_{k+1} \in A \mid \tilde{Z}_k = x)$ has density

$$\tilde{p}_k(x, y) = \frac{p(y, x) q_{\tau - k - 1}(y)}{q_{\tau - k}(x)}$$

with respect to $\mu$, where $q_n$ is the p.d.f. of $Q_n = Q_{\text{true}} P^n$, the distribution of $Z_n$.

The distributions $Q_n, n \geq 0$ are unknown, since they depend on the data distribution $P_{\text{true}}$, and the transition probabilities above must be estimated from data to provide a sampling algorithm from the reversed Markov chain. While, at first glance, this does not seem like a simplification of the problem, because one now has to sample from a potentially large number ($\tau$) of distributions instead of one, this leads, with proper modeling and some intensive learning, to efficient and accurate sampling algorithms.

Several factors can indeed make this approach achievable. First, the forward chain should be making small changes to the current configuration at each step (e.g., adding a small amount of noise). This ensures that the reversed transition probabilities $\tilde{p}_k(x, \cdot)$ are close to Dirac distributions and are therefore likely to be well approximated by simple unimodal distributions such as Gaussians. Second, the estimation problem does not have hidden data: given an observed sample, one can simulate $\tau$ steps of the forward chain to obtain, after reversing the order, a full observation of the reversed chain. Third, in some cases, analytical considerations can lead to partial computations that facilitate the modeling of the reversed transitions.

### 18.4.2 Binary model

We now take some examples, starting with a discrete one. Let $Q_{\text{true}}$ be the distribution of a binary random field with state space $\{0, 1\}$ over a set of vertexes $V$, i.e., with the notation of section 13.2, $\mathcal{R}_X = \mathcal{F}(V)$ with $F = \{0, 1\}$. Fix a small $\epsilon > 0$ and define the transition probability $p(x, y)$ for $x, y \in \mathcal{F}(V)$ by

$$p(x, y) = \prod_{s \in V} \left( (1 - \epsilon) \mathbf{1}_{y^{(s)} = x^{(s)}} + \epsilon \mathbf{1}_{y^{(s)} = 1 - x^{(s)}} \right).$$

Since $p(x, y) > 0$ for all $x$ and $y$, the chain converges (uniformly geometrically) to its invariant probability $Q_\infty$ and one easily checks that this probability is such that all variables are independent Bernoulli random variables with success probability $1/2$. Assuming that $\tau$ is large enough so that $Q_\tau \simeq Q_\infty$, the sampling algorithm initializes the reversed chain as independent Bernoulli($1/2$) variables and runs $\tau$ steps using the transitions $\tilde{p}_k$ which must be learned from data.

For this model, we have

$$q_k(x) = \sum_{y \in \mathcal{F}(V)} q_{k-1}(y) p(y, x).$$

For this transition, the probabililty of flipping two or more values of $y$ is

$$1 - (1 - \epsilon)^N - N\epsilon(1 - \epsilon)^{N-1} = \frac{N(N-1)}{2}\epsilon^2 + o(\epsilon^2)$$

with $N = |V|$. We will write $x \sim_s y$ if $y^{(s)} = 1 - x^{(s)}$ and $y^{(t)} = x^{(t)}$ for $s \neq t$, and we will write $x \sim y$ if $x \sim_s y$ for some $s$. With this notation, we have

$$q_k(x) = (1 - N\epsilon)q_{k-1}(x) + \epsilon \sum_{y:y \sim x} q_{k-1}(y) + O(\epsilon^2)$$

Since it implies that $q_k(x) = q_{k-1}(x) + o(\epsilon)$, this expression can be reversed as

$$q_{k-1}(y) = (1 + N\epsilon)q_k(y) - \epsilon \sum_{x:x \sim y} q_k(x) + O(\epsilon^2)$$

Similarly, we have

$$p(y, x) = (1 - N\epsilon)\mathbf{1}_{x=y} + \epsilon\mathbf{1}_{x \sim y} + O(\epsilon^2).$$

This gives

$$p(y, x)q_{k-1}(y) = q_k(x)\mathbf{1}_{x=y} - \epsilon\mathbf{1}_{x=y} \sum_{x':x' \sim y} q_k(x') + \epsilon q_k(y)\mathbf{1}_{x \sim y} + O(\epsilon^2),$$

and we finally get

$$\tilde{p}_k(x, y) = \left(1 - \epsilon \sum_{x':x' \sim y} \frac{q_{\tau-k}(x')}{q_{\tau-k}(x)}\right)\mathbf{1}_{x=y} + \epsilon\frac{q_{\tau-k}(y)}{q_{\tau-k}(x)}\mathbf{1}_{x \sim y} + O(\epsilon^2)$$

If one lets $\sigma_k^{(s)}(x) = \frac{q_{\tau-k}(y)}{q_{\tau-k}(x)}$ with $y \sim_s x$, and defines

$$\hat{p}_k(x, y) = \prod_{s \in V}\left((1 - \epsilon\sigma_k^{(s)}(x))\mathbf{1}_{y^{(s)}=x^{(s)}} + \epsilon\sigma_k^{(s)}(x)\mathbf{1}_{y^{(s)}=1-x^{(s)}}\right),$$

one checks easily that $\hat{p}_k(x, y) = \tilde{p}_k(x, y) + O(\epsilon^2)$. This suggests modeling the reversed chain using transitions $\hat{p}_k$, for which the mapping $x \mapsto (\sigma_k^{(s)}(x), s \in V)$ needs to be learned from data (for example using a deep neural network). Note that $1 - \sigma_k(x)$ is precisely the score function introduced for discrete distributions in remark 17.8.

### 18.4.3   Model with continuous variables

We now switch to an example with vector-valued variables, $\mathcal{R}_X = \mathbb{R}^d$, and assume that the forward Markov chain is such that, conditionally to $X_n = x$,

$$X_{n+1} \sim \mathcal{N}(x + hf(x), \sqrt{h}\mathrm{Id}_{\mathbb{R}^d}),$$

where $f$ is $C^1$. We saw in section 12.3.7 that, when $f = -\nabla H/2$ for a $C^2$ function $H$ such that $\exp(-H)$ is integrable, this chain converges (approximately for small $h$) to a limit distribution with p.d.f. (with respect to Lebesgue's measure) proportional to $\exp(-H)$. In the linear case, in which $f(x) = -Ax/2$ for some positive-definite symmetric matrix $A$, so that $H(x) = \frac{1}{2}x^T A x$, the limit distribution can be identified exactly as $\mathcal{N}(0, \Sigma_h)$ where $\Sigma_h$ satisfies the equation

$$A\Sigma_h + \Sigma_h A - \frac{h}{2}A^2 - 2\mathrm{Id}_{\mathbb{R}^d} = 0$$

whose solution is $\Sigma_h = (A - hA^2/4)^{-1}$ (details being left to the reader). This implies that this limit distribution can be easily sampled from for any choice of $A$.

We now return to general $f$'s and make, like in the discrete case, a first-order identification of the reversed chain. We note that, for any smooth function $\gamma$,

$$\mathbb{E}(\gamma(X_{n+1}) \mid X_n = x) = \mathbb{E}(\gamma(x + hf(x) + \sqrt{h}U))$$

where $U \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$. Making the second order expansion

$$\gamma(x + hf(x) + hU) = \gamma(x) + \sqrt{h}\nabla\gamma(x)^T U + h\nabla\gamma(x)^T f(x) + \frac{h}{2}U^T \nabla^2\gamma(x)U + o(h)$$

and taking the expectation gives

$$\mathbb{E}(\gamma(X_{n+1}) \mid X_n = x) = \gamma(x) + h\nabla\gamma(x)^T f(x) + \frac{h}{2}\Delta\gamma(x) + o(h). \tag{18.10}$$

Considering the reversed chain, and letting $q_k$ denote the p.d.f. of $X_k$ for the forward chain, we have

$$\begin{aligned}
\mathbb{E}(\gamma(X_{k-1}) \mid X_k = x) &= \int_{\mathbb{R}^d} \gamma(y)\tilde{p}_k(x,y)dy \\
&= \int_{\mathbb{R}^d} \gamma(y)p(y,x)\frac{q_{k-1}(y)}{q_k(x)}dy \\
&= \frac{1}{(2\pi h)^{d/2}} \int_{\mathbb{R}^d} \gamma(y)\frac{q_{k-1}(y)}{q_k(x)}e^{-\frac{1}{2h}|x-y-hf(y)|^2}dy \\
&= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} \gamma(x - \sqrt{h}u)\frac{q_{k-1}(x - \sqrt{h}u)}{q_k(x)}e^{-\frac{1}{2}|u - \sqrt{h}f(x - \sqrt{h}u)|^2}dy,
\end{aligned}$$

with the change of variable $u = (x - y)/\sqrt{h}$. We make a first-order expansion of the terms in this integral, with

$$\gamma(x - \sqrt{h}u)q_{k-1}(x - \sqrt{h}u) = \gamma(x)q_{k-1}(x) - \sqrt{h}\nabla(\gamma q_{k-1})(x)^T u + \frac{h}{2}u^T\nabla^2(\gamma q_{k-1})(x)u + o(h)$$

and

$$e^{-\frac{1}{2}|u - \sqrt{h}f(x - \sqrt{h}u)|^2} = e^{-\frac{1}{2}|u|^2}e^{\sqrt{h}u^T f(x) - hu^T df(x)u - \frac{1}{2}|f(x)|^2 + o(h)}$$

$$= e^{-\frac{1}{2}|u|^2}\left(1 + \sqrt{h}u^T f(x) - hu^T df(x)u - \frac{h}{2}|f(x)|^2 + \frac{h}{2}|u^T f(x)|^2 + o(h)\right).$$

Taking products

$$\gamma(x - \sqrt{h}u)q_{k-1}(x - \sqrt{h}u)e^{-\frac{1}{2}|u - \sqrt{h}f(x - \sqrt{h}u)|^2}$$

$$= e^{-\frac{1}{2}|u|^2}\gamma(x)q_{k-1}(x)\left(1 + \sqrt{h}u^T f(x) - hu^T df(x)u - \frac{h}{2}|f(x)|^2 + \frac{h}{2}|u^T f(x)|^2\right)$$

$$+ e^{-\frac{1}{2}|u|^2}\left(-\sqrt{h}\nabla(\gamma q_{k-1})(x)^T u - h(\nabla(\gamma q_{k-1})(x)^T u)(f(x)^T u) + \frac{h}{2}u^T\nabla^2(\gamma q_{k-1})(x)u\right) + o(h)$$

We now take the integral with respect to $u$ (recall that $\mathbb{E}(U^T A U) = \text{trace}(A)$ if $A$ is any square matrix and $U$ is standard Gaussian), so that

$$\frac{1}{(2\pi)^{d/2}}\int_{\mathbb{R}^d}\gamma(x - \sqrt{h}u)q_{k-1}(x - \sqrt{h}u)e^{-\frac{1}{2}|u - \sqrt{h}f(x - \sqrt{h}u)|^2}\,du$$

$$= \gamma(x)q_{k-1}(x) + h\left(-\gamma(x)q_{k-1}(x)\nabla\cdot f(x) - \nabla(\gamma q_{k-1})(x)^T f(x) + \frac{1}{2}\Delta(\gamma q_{k-1})(x)\right) + o(h)$$

$$= q_{k-1}(x)\left(\gamma(x) + h\left(-\gamma(x)\nabla\cdot f(x) - \left(\frac{\nabla(\gamma q_{k-1})(x)}{q_{k-1}(x)}\right)^T f(x) + \frac{1}{2}\frac{\Delta(\gamma q_{k-1})(x)}{q_{k-1}(x)}\right)\right) + o(h)$$

To compute an expansion of $q_k(x)$, it suffices to take $\gamma = 1$ above, so that

$$q_k(x) = q_{k-1}(x)\left(1 + h\left(-\nabla\cdot f(x) - \left(\frac{\nabla q_{k-1}(x)}{q_{k-1}(x)}\right)^T f(x) + \frac{1}{2}\frac{\Delta q_{k-1}(x)}{q_{k-1}(x)}\right)\right) + o(h).$$

We now take the first-order expansion of the ratio, removing terms that cancel, and get

$$\mathbb{E}(\gamma(X_{k-1})\,|\,X_k = x) = \gamma(x) - h\nabla\gamma(x)^T f(x) + h\nabla\gamma(x)^T\left(\frac{\nabla q_{k-1}(x)}{q_{k-1}(x)}\right) + \frac{h}{2}\Delta\gamma(x) + o(h)$$

Comparing with (18.10), we find that $\tilde{X}_k = X_{\tau-k}$ behaves, for small $h$, like the non-homogeneous Markov chain such that the conditional distribution of $\tilde{X}_{k+1}$ given $\tilde{X}_k = x$ is $\mathcal{N}(x - hf(x) - hs_{\tau-k-1}(x), \sqrt{h}\text{Id}_{\mathbb{R}^d})$, with $s_{\tau-k-1}(x) = -\nabla\log q_{\tau-k-1}$, the score function introduced in section 17.2.6, and score-matching methods from that section can be used to estimate it from observations of the forward chain initialized with training data.

### 18.4.4 Continuous-time limit

The forward schemes described in the previous examples can be interpreted as continuous time processes over discrete or continuous variables. In the latter case, the example $X_{k+1} \sim \mathcal{N}(x + hf(x), \sqrt{h}\mathrm{Id}_{\mathbb{R}^d})$ conditionally to $X_k = x$ is a discretization of the stochastic differential equation

$$dx_t = f(x_t)dt + dw_t$$

(see remark 12.5), where $w_t$ is a Brownian motion and the diffusion is initialized with $Q_{\mathrm{true}}$. We found that going backward meant (at first order and conditionally to $X_k = x$)

$$X_{k-1} \sim \mathcal{N}(x - hf(x) - hs_{k-1}(x), \sqrt{h}\mathrm{Id})$$

that we can rewrite as

$$x_\tau - X_{k-1} \sim \mathcal{N}(x_\tau - x + hf(x) + hs_{k-1}(x), \sqrt{h}\mathrm{Id}).$$

Following the definition in Anderson [9], this corresponds to a first-order discretization of the *reverse diffusion*

$$dx_t = (f(x_t) + s_t(x_t))dt + d\tilde{w}_t, \ t \leq \tau$$

where $\tilde{w}_t$ is also a Brownian motion. This reverse diffusion with $X_\tau \sim Q_\infty$ will therefore approximately sample from $Q_{\mathrm{true}}$. (With this terminology, forward and reverse diffusions have similar differential notation, but mean different things.) Note that, in the continuous-time limit, the reverse Markov process follows the distribution of the reversed diffusion exactly.

### 18.4.5 Differential of neural functions

As we have seen in the previous two examples, estimating the reversed Markov chain requires computing the score functions of the forward probabilities. In the case of continuous variables, this score function is typically parametrized as a neural network, so that the function $s_k(x) = -\nabla \log q_k(x)$ is computed as $s_k(x) = F(x; \mathcal{W}_k)$, with the usual definition $F(x, \mathcal{W}_k) = z_{m+1}$ with $z_{j+1} = \varphi_j(z_j, w_{jk})$, $z_0 = x$ and $\mathcal{W}_k = (w_{0k}, \dots, w_{mk})$.

Assume that a training set $T$ is observed. Running the forward Markov chain initialized with elements of $T$ generates a new training step at each time step, that we will denote $T_k$ at step $k$. We have seen in section 17.2.6 that the score function $s_k$ could be estimated by minimizing, with respect to $\mathcal{W}$

$$\sum_{x \in T_k} \left( |F(x, \mathcal{W})|^2 - 2\nabla \cdot F(x, \mathcal{W}) \right).$$

This term involves the differential of $F$, which is defined recursively by (simply taking the derivative at each step)

$$dF(x, \mathcal{W}) = \zeta_{m+1}, \quad \zeta_{j+1} = d\varphi_j(z_j, w_j)\zeta_j,$$

with $\zeta_0 = \mathrm{Id}_{\mathbb{R}^d}$. From this recursive definition, back-propagation can be applied, in principle, to compute the derivative of $dF(x, \mathcal{W})$ with respect to $\mathcal{W}$. The feasibility of this computation, however, is limited when $d$ is large ($d$ could be tens of thousands if one models images) computing the $d \times d$ matrix $dF(x, \mathcal{W})$ is intractable.

We can note that, for any $h \in \mathbb{R}^d$, the vector $dF(x, \mathcal{W})h$ also satisfies the recursion

$$dF(x, \mathcal{W})h = \zeta_{m+1}h, \quad \zeta_{j+1}h = d\varphi_j(z_j, w_j)\zeta_j h,$$

with $\zeta_0 h = h$ and

$$\nabla \cdot F(x, \mathcal{W}) = \sum_{i=1}^{d} \mathfrak{e}_i^T \, dF(x, \mathcal{W}) \mathfrak{e}_i$$

where $\mathfrak{e}_1, \ldots, \mathfrak{e}_d$ is the canonical basis of $\mathbb{R}^d$. Putting the divergence of $F$ in this form does not reduce the computation cost (which is, roughly $d^2 m$, assuming that all $z_j$'s have the same dimension), but expresses the divergence term in a form that is amenable to stochastic gradient descent (which is typically already used to approximate the sum over $x$). Indeed, if $U$ follows any distribution with zero mean and covariance matrix equal to the identity (such as a standard Gaussian, or the uniform distribution on the unit sphere), then

$$\nabla \cdot F(x, \mathcal{W}) = \mathbb{E}(U^T dF(x, \mathcal{W})U)$$

so that $U$ can be sampled from in minibatches in SGD implementations (see [180], where this approach is called "sliced score matching").

# Chapter 19

# Clustering

## 19.1 Introduction

We now describe a collection of methods designed to divide a training set into homogeneous subsets, or clusters. This grouping operation is a key problem in many applications for which it is important to categorize the data in order to obtain improved understanding of the sampled phenomenon, and sometimes to be able to apply a different approach to subsequent processing or analysis adapted to each cluster.

We will assume that the variables of interest belong a set $\mathcal{R} = \mathcal{R}_X$ where $\mathcal{R}$ is equipped with a discrepancy function $\alpha : \mathcal{R} \times \mathcal{R} \to [0, +\infty)$. Often, $\alpha$ is derived from a distance $\rho$ on $\mathcal{R}$, but this is not always the case. We will assume that the data results from a training set $T = (x_1, \ldots, x_N)$. However, it may happen that only the discrepancy matrix $A = (\alpha(x, y), x, y \in T)$ is observed, while a coordinate representation of the elements of $T$ is not available.

Let us consider a few examples.

(i) The simplest case is when $\mathcal{R} = \mathbb{R}^d$ with the standard Euclidean metric. Slightly more generally, a metric may be defined by $\rho^2(x, y) = \|h(x) - h(y)\|_H^2$, where $H$ is an inner-product space and the feature function $h : \mathcal{R} \mapsto H$ may be unknown, while its associated "kernel", $K(x, y) = \langle h(x), h(y) \rangle_H$ is known (this is a metric if $h$ is one-to-one). In this case

$$\rho^2(x, y) = K(x, x) - 2K(x, y) + K(y, y).$$

Typically, one then takes $\alpha = \rho$ or $\alpha = \rho^2$.

(ii) Very often, however, the data is not Euclidean, and the distance does not correspond to a feature space representation. This is the case, for example, for data belonging to "curved spaces" (manifolds), for which one may use the intrinsic distance

451

provided by the length of shortest paths linking two points (assuming of course that this notion can be given a rigorous meaning). The simplest example is data on the unit sphere, where the distance $\rho(x,y)$ between two points $x$ and $y$ is the length of the shortest large circle that connects them, satisfying

$$|x - y|^2 = 2 - 2\cos\rho(x,y).$$

Once again, $\alpha = \rho$ or $\rho^2$ is a typical choice.

(iii) A more complex example is provided by $\mathcal{R}$ being the space of symmetric positive-definite matrices on $\mathbb{R}^d$, for which one defines the length of a differentiable curve $(S(t), t \in [a,b])$ in this space by

$$\int_a^b \sqrt{\text{trace}((S(t)^{-1}\partial_t S)(S(t)^{-1}\partial_t S)^T)}\,dt$$

and for which

$$\rho^2(S_1, S_2) = \sum_{i=1}^d (\log \lambda_i)^2$$

where $\lambda_1, \ldots, \lambda_d$ are the eigenvalues of $S_1^{-1/2} S_2 S_1^{-1/2}$ or, equivalently, solutions of the generalized eigenvalue problem $S_2 u = \lambda S_1 u$ (see, for example, [72]).

(iv) Another common assumption is that the elements of $\mathcal{R}$ are vertices of a weighted graph of which $T$ is a subgraph; $\rho$ may then be, e.g., the geodesic distance on the graph.

## 19.2  Hierarchical clustering and dendograms

### 19.2.1  Partition trees

This method builds clusters by organizing them in a binary hierarchy in which the data is divided into subsets, starting with the full training set, and iteratively splitting each subset into two parts until reaching singletons. This results in a binary tree structure, called a *dendogram*, or partition tree, which is defined as follows.

**Definition 19.1** *A partition tree of a finite set A is a finite collection of nodes $\mathcal{T}$ with the following properties.*

*(i) Each node has either zero or exactly two children. (We will use the notation $v \to v'$ to indicate that $v'$ is a child of $v$.)*

*(ii) All nodes but one have exactly one parent. The node without parent is the* root *of the tree.*

*(iii) To each node $v \in \mathcal{T}$ is associated a subset $A_v \subset A$.*

1: $\{a, b, c, d, e, f\}$

2: $\{a, c, f\}$

3: $\{b, d, e\}$

4: $\{a, f\}$

5: $\{c\}$

6: $\{d\}$

7: $\{b, e\}$

8: $\{a\}$

9: $\{f\}$

10: $\{b\}$

11: $\{e\}$

Figure 19.1: A partition tree of the set $\{a, b, c, d, e, f\}$.

*(iv) If $v'$ and $v''$ are the children of $v$, then $(A_{v'}, A_{v''})$ forms a partition of $A_v$.*

*Nodes without children are called leaves, or terminal nodes. We will say that the hierarchy is complete if $A_v = A$ if $v$ is the root, and $|A_v| = 1$ for all terminal nodes.*

An example of partition tree is provided in fig. 19.1.

The construction of the tree can follow two directions, the first one being bottom-up, or agglomerative, in which the algorithm starts with the collection of all singletons and merges subsets one pair at a time until everything is merged into the full dataset. The second approach is top-down, or divisive, and initializes the algorithm with the full training set which is recursively split until singletons are reached. The first approach, on which we now focus, is more common, and computationally simpler.

We let $T$ denote the training set and assume that a matrix of dissimilarities

$$(\alpha(x, y), x, y \in T)$$

is given. We will make the abuse of notation of considering that $T$ is a set even though some of its elements may be repeated. This is no loss of generality, since $T = (x_1, \ldots, x_N)$ can always be replaced by the subset $\{(k, x_k), k = 1, \ldots, N\}$ of $\mathbb{N} \times \mathcal{R}$.

### 19.2.2 Bottom-up construction

We will extend $\alpha$ to a dissimilarity measure between subsets $A, A' \subset T$ that we will denote $(A, A') \mapsto \varphi(A, A')$. Once $\varphi$ is defined, agglomeration works along the following algorithm.

---

**Algorithm 19.1**

1. Start with the collection $\mathcal{T}_1, \ldots, \mathcal{T}_N$ of all single-node trees associated to each element of $T$. Let $n = 0$ and $m = N$.

2. Assume that, at step $n$ of the algorithm, one has a collection of partition trees $\mathcal{T}_1,\ldots,\mathcal{T}_m$ with root nodes $r_1,\ldots,r_m$ associated with subsets $A_{r_1},\ldots,A_{r_m}$ of $T$. Let the total collection of nodes be indexed as $\mathcal{V}_n = \{v_1,\ldots,v_{N+n}\}$, so that $\{r_1,\ldots,r_m\} \subset \mathcal{V}_n$.

3. If $m = 1$, stop the algorithm.

4. Select indices $i,j \in \{1,\ldots,m\}$ such that $\varphi(A_{r_i}, A_{r_j})$ is minimal, and merge the corresponding trees by creating a new node $v_{n+1+N}$ with the root nodes of $\mathcal{T}_i$ and $\mathcal{T}_j$ as children (so that $v_{n+1+N}$ is associated with $A_{r_i} \cup A_{r_j}$). Add $v_{n+1+N}$ to the collection of root nodes, and remove $r_i$ and $r_j$.

5. Set $n \to n+1$ and $m \to m-1$ and return to step 2.

---

Clearly, the specification of the extended dissimilarity measure ($\varphi$) is a key element of the method. Some of most commonly used extensions are:

- Minimum gap: $\varphi_{\min}(A, A') = \min(\alpha(x, x') : x \in A, x' \in A')$.

- Maximum dissimilarity: $\varphi_{\max}(A, A') = \max(\alpha(x, x') : x \in A, x' \in A')$.

- Sum of dissimilarities:

$$\varphi_{\text{sum}}(A, A') = \sum_{x \in A} \sum_{x' \in A'} \alpha(x, x')$$

- Average dissimilarity:

$$\varphi_{\text{avg}}(A, A') = \frac{1}{|A||A'|} \sum_{x \in A} \sum_{x' \in A'} \alpha(x, x').$$

As shown in the next two propositions, the maximum distance favors clusters with small diameters, while using minimum gaps tends to favor connected clusters.

**Proposition 19.2** *Let* $\text{diam}(A) = \max(\alpha(x, y), x, y \in A)$. *The agglomerative algorithm using* $\varphi_{\max}$ *is identical to that using* $\varphi(A, A') = \text{diam}(A \cup A')$.

PROOF Call Algorithm 1 the agglomerative algorithm using $\varphi_{\max}$, and Algorithm 2 the one using $\varphi$. At initialization, we have (because all sets are singletons),

$$\varphi_{\max}(A_k, A_l) = \text{diam}(A_k \cup A_l) \text{ for all } 1 \le k \ne l \le m. \tag{19.1}$$

We show that this property remains true at all steps of the algorithms. Proceeding by induction, assume that, up to the step $n$, Algorithms 1 and 2 have been identical and result in sets $(A_1,\ldots,A_m)$ satisfy (19.1). Then the next steps of the two algorithms coincide and assume, without loss of generality, that this next step

merges $A_{m-1}$ with $A_m$. Let $A'_{m-1} = A_{m-1} \cup A_m$ so that $\mathrm{diam}(A'_{m-1}) \leq \mathrm{diam}(A_i \cup A_j)$ for all $1 \leq i \neq j \leq m$.

We need to show that the new partition satisfies (19.1), which requires that

$$\varphi_{\max}(A'_{m-1}, A_k) = \mathrm{diam}(A'_{m-1} \cup A_k)$$

for $k = 1, \ldots, m - 2$.

We have

$$\mathrm{diam}(A'_{m-1} \cup A_k) = \max(\mathrm{diam}(A'_{m-1}), \mathrm{diam}(A_k), \varphi_{\max}(A'_{m-1}, A_k)),$$

so that we must show that

$$\max(\mathrm{diam}(A'_{m-1}), \mathrm{diam}(A_k)) \leq \varphi_{\max}(A'_{m-1}, A_k).$$

Write

$$\varphi_{\max}(A'_{m-1}, A_k) = \max(\varphi_{\max}(A_m, A_k), \varphi_{\max}(A_{m-1}, A_k))$$
$$= \max(\mathrm{diam}(A_m \cup A_k), \mathrm{diam}(A_{m-1} \cup A_k))$$

where the last identity results from the induction hypothesis.

The fact that

$$\mathrm{diam}(A_k) \leq \max(\mathrm{diam}(A_m \cup A_k), \mathrm{diam}(A_{m-1} \cup A_k))$$

is obvious, and the inequality

$$\mathrm{diam}(A'_{m-1}) \leq \max(\mathrm{diam}(A_m \cup A_k), \mathrm{diam}(A_{m-1} \cup A_k))$$

results from the fact that $A_m$ and $A_{m-1}$ was an optimal pair. This shows that the induction hypothesis remains true at the next step and concludes the proof of the proposition. ∎

We now analyze $\varphi_{\min}$ and, more specifically, the equivalence between the resulting algorithm and the one using the following measure of connectedness. For a given set $A$ and $x, y \in A$, let

$$\tilde{\alpha}_A(x, y) = \inf \left\{ \epsilon : \exists n > 0, \exists (x = x_0, x_1, \ldots, x_{n-1}, x_n = y) \in A^{n+1} : \right.$$

$$\left. \alpha(x_i, x_{i-1}) \leq \epsilon \text{ for } 1 \leq i \leq n \right\}.$$

So $\tilde{\alpha}_A$ is the smallest $\epsilon$ such that there exists a sequence of steps of size less than $\epsilon$ in $A$ going from $x$ to $y$. The function

$$\mathrm{conn}(A) = \max\{\tilde{\alpha}_A(x, y) : x, y \in A\}$$

measures how well the set $A$ is connected relative to the dissimilarity measure $\alpha$. and we have:

**Proposition 19.3** *The agglomerative algorithm using* $\varphi_{\min}$ *is identical to that using* $\varphi(A, A') = \text{conn}(A \cup A')$.

PROOF  The proof is similar to that of proposition 19.2. Indeed one can note that

$$\text{conn}(A \cup A') = \max(\text{conn}(A), \text{conn}(A'), \varphi_{\min}(A, A')).$$

Given this we can proceed by induction and prove that, if the current decomposition is $A_1, \ldots, A_m$ such that $\psi(A_k \cup A_l) = \varphi_{\min}(A_k, A_l)$ for all $1 \le k \ne l \le m$, then this property is still true after merging using $\varphi_{\min}$ and $\varphi$.

Assuming again that $A_{m-1}$ and $A_m$ are merged, and letting $A'_{m-1} = A_m \cup A_{m-1}$, we need to show that $\text{conn}(A_k \cup A'_{m-1}) = \varphi_{\min}(A_k, A'_{m-1})$ for all $k = 1, \ldots, m-2$, which is the same as showing that:

$$\max(\text{conn}(A_k), \text{conn}(A'_{m-1})) \le \varphi_{\min}(A_k, A'_{m-1}) = \min(\varphi_{\min}(A_k, A_{m-1}), \varphi_{\min}(A_k, A_m)).$$

From the induction hypothesis, we have

$$\min(\varphi_{\min}(A_k, A_{m-1}), \varphi_{\min}(A_k, A_m)) = \min(\text{conn}(A_k \cup A_{m-1}), \text{conn}(A_k \cup A_m))$$

and both terms in the right-hand side are larger than $\text{conn}(A_k)$ and also larger than $\text{conn}(A'_{m-1})$ which was a minimizer.                                                                    ∎

### 19.2.3   Top-down construction

The agglomerative method is the most common way to build dendograms, mostly because of the simplicity of the construction algorithm. The divisive approach is more complex, because the division step, which requires, given a set $A$, to optimize a splitting criterion over all two-partitions of $A$, may be significantly more expensive than the merging steps in the agglomerative algorithm. The top-down construction therefore requires the specification of a "splitting algorithm" $\sigma : A \mapsto (A', A'')$ such that $(A', A'')$ is a partition of $A$. We assume that, if $|A| > 1$, then the partition $A, A''$ is not trivial, i.e., neither set is empty.

Given $\sigma$, the top-down construction is as follows.

---

**Algorithm 19.2**

   1. Start with the one-node partition tree $\mathcal{T}_0 = (T)$.

   2. Assume that at a given step of the algorithm, the current partition is $\mathcal{T}$.

   3. If $\mathcal{T}$ is complete, stop the algorithm.

   4. For each terminal node $v$ in $\mathcal{T}$ such that $|A_v| > 1$, compute $(A'_v, A''_v) = \sigma(A_v)$ and add two children $v'$ and $v''$ to $v$ with $A_{v'} = A'_v$ and $A_{v''} = A''_v$.

   5. Return to step 2.

---

The division of a set into two parts is itself a clustering algorithm, and one may apply any of those described in the rest of this chapter.

### 19.2.4 Thresholding

Once a complete hierarchy is built, it provides a complete binary partition tree $\mathcal{T}$. This tree provides in turn a collection of partitions of $\mathcal{V}$, each of them obtained through pruning. We now formalize this operation.

Let $\mathcal{V}_T$ denote the set of terminal nodes in $\mathcal{T}$ and $\mathcal{V}_0 = \mathcal{V} \setminus \mathcal{V}_T$ contain the interior nodes. Define a *pruning set* to be a subset $\mathcal{D} \subset \mathcal{V}_0$ that contains no pair of nodes $v, v'$ such that $v'$ is a descendant of $v$. To any pruning set $\mathcal{D}$, one can associate the pruned subtree $\mathcal{T}(\mathcal{D})$ of $\mathcal{T}$ consisting of $\mathcal{T}$ from which all the vertices that are descendants of elements of $\mathcal{D}$ are removed. From any such pruned subtree, one obtain a partition $S(\mathcal{D})$ of $T$ formed by the collection of sets $A_v$ for $v$ in the terminal nodes of $\mathcal{T}(\mathcal{D})$. Between the extreme case $S(v_0) = \{\mathcal{V}\}$ (where $v_0$ is the root of $\mathcal{T}$) and $S(\emptyset) = (\{x\}, x \in \mathcal{V}_T)$, there exists a huge number of possible partitions obtained in this way.

It is often convenient to organize these partitions according to the level sets of a well-chosen score function $v \mapsto h(v)$ defined over $\mathcal{V}_0$. For $\mathcal{D} \subset \mathcal{V}$, we denote by $\max(\mathcal{D})$ the set of its deepest elements, i.e., the set formed by those $v \in \mathcal{D}$ that have no descendant in $\mathcal{D}$. Then, for any $\lambda \in \mathbb{R}$, one can define $\mathcal{D}_\lambda^+ = \max\{v : h(v) \geq \lambda\}$ (resp. $\mathcal{D}_\lambda^- = \max\{v : h(v) \leq \lambda\}$) and the associated partition $S(\mathcal{D}_\lambda^+)$ (resp. $S(\mathcal{D}_\lambda^-)$). The score function $h$ can be linked to the construction algorithm. For example, if one uses a bottom-up construction using an extended dissimilarity $\varphi$, one can associate to each node $v$ with $v \in \mathcal{V}_0$ the value of $\varphi(A_{v'}, A_{v''})$ where $v'$ and $v''$ are the children of $v$.

Another way to define such scores functions is by assigning weights to edges in $\mathcal{T}$. Indeed, given a collection $w$ of positive numbers $w(v, v')$ for $v \to v'$ in $\mathcal{T}$, one can define a score $h_w$ recursively by letting $h_w(v_0) = 0$ and $h_w(v') = h_w(v) + w(v, v')$ if $v'$ is a child of $v$. The choice $w(v, v') = 1$ for all $v, v'$ provide the usual notion of depth in the tree.

Scores can also be built bottom-up, letting $h(v) = 0$ for terminal nodes and, for $v \in \mathcal{V}_0$,

$$h_w(v) = \max(h_w(v') + w(v, v'), h_w(v'') + w(v, v''))$$

where $v', v''$ are the children of $v$ Here, taking $w = 1$ provides the height of each node.

## 19.3 K-medoids and K-mean

### 19.3.1 K-medoids

One of the limitation of hierarchical clustering is that it is a greedy approach that does not optimize a global quality measure associated to the partition. Such qual-

ity measures can indeed be defined based on the heuristic that clusters should be homogeneous (for some criterion) and far apart from each other.

In centroid-based methods, the homogeneity criterion is the minimum, over all possible points in $\mathcal{R}$, of the sum of dissimilarities between elements of the cluster and that point. More precisely, for any $A \subset \mathcal{R}$, and any dissimilarity measure $\alpha$, define the *central dispersion* index

$$V_\alpha(A) = \inf \left\{ \sum_{x \in A} \alpha(x, c) : c \in \mathcal{R} \right\}. \tag{19.2}$$

If $c$ achieves the minimum in the definition of $V_\alpha$, it is called a *centroid* of $A$ for the dissimilarity $\alpha$.

The most common choice is $\alpha = \rho^2$, where $\rho$ is a metric on $\mathcal{R}$, and in this case, we will just use $V$ in place of $V_{\rho^2}$. Note also that it is always possible to limit $\mathcal{R}$ to the training set $T$, in which case the optimization in (19.2) is over a finite number of centers. This makes centroid-based methods also applicable to the situation when the matrix of dissimilarities is the only input provided to the algorithm, or when the set $\mathcal{R}$ and the function $\alpha$ are too complex for the optimization in (19.2) to be feasible.

A centroid, $c$, in (19.2) may not always exists, and when it exists it may not always be unique. For $\alpha = \rho^2$, a point $c$ such that

$$V(A) = \sum_{x \in A} \rho^2(x, c)$$

is called a Fréchet mean of the set $A$. Returning to the examples provided in the beginning of this chapter, two antipodal points on the sphere (whose distance is $\pi$) have an infinity of Fréchet means (or midpoints in this case) provided by every point in the equator between them. In contrast, the example provided with symmetric matrices provides a so-called Hadamard space [44] and the Fréchet mean in that case is unique. Of course, for Euclidean metrics, the Fréchet mean is just the usual one.

Returning to our general discussion, the K-medoids method optimizes the sum of central dispersions with a fixed number of clusters. Note that the letter K in K-medoids originally refers to this number of clusters, but this notation conflicts with other notation in this book (e.g., reproducing kernels) and we shall denote by $p$ this

target number[1]. So the K-medoids method minimizes

$$W_\alpha(A_1,\ldots,A_p) = \sum_{i=1}^{p} V_\alpha(A_i)$$

over all partitions $A_1,\ldots,A_p$ of the training set $T$. Equivalently, it minimizes

$$\mathbb{W}_\alpha(A_1,\ldots,A_p,c_1,\ldots,c_p) = \sum_{i=1}^{p} \sum_{x \in A_i} \alpha(x,c_i) \qquad (19.3)$$

over all partitions of $T$ and $c_1,\ldots,c_p \in \mathcal{R}$. Finally, taking first the minimum with respect to $A_i$, which corresponds to associating each $x$ to the subset with closest center, K-medoids, an equivalent formulation minimizes

$$\tilde{W}_\alpha(c_1,\ldots,c_p) = \sum_{x \in T} \min\left\{\alpha(x,c_i), i = 1,\ldots,p\right\}.$$

The standard implementation of K-medoids solves this problem using an alternate minimization, as defined in the following algorithm.

---

**Algorithm 19.3 (K-medoids)**

Let $T \subset \mathcal{R}$ be the training set. Start with an initial choice of $c_1,\ldots,c_p \in \mathcal{R}$ and iterate over the following two steps until stabilization:

(1) For $i = 1,\ldots,p$, let $A_i$ contain points $x \in T$ such that $\alpha(x,c_i) = \min\{\alpha(x,c_j), j = 1,\ldots,p\}$. In case of a tie in this minimum, assign $x$ to only one of the tied sets (e.g., at random) to ensure that $A_1,\ldots,A_p$ is a partition.

(2) For $i = 1,\ldots,p$, let $c_i$ be a minimizer of $\sum_{x \in A_i} \alpha(x,c_i)$ if $A_i$ is not empty, or $c_i$ be a random point in $T$ otherwise.

---

It should be clear that each step reduces the total cost $\mathbb{W}_\alpha$ and that this cost should stabilize at some point (which provides the stopping criterion) because there is only a finite number of possible partitions of $T$. However, there can be many possible limit points that are stable under the previous iterations, and some may correspond to poor "local minima" of the objective function. Since the end-point of the algorithm depends on the initialization, this step requires extra care. One may design *ad-hoc* heuristics in order to start the algorithm with a good initial point that is likely to provide a good solution at the end. These heuristics may depend on the

---

[1]We still call the method K-medoids rather than $p$-medoids, to keep the name universally used in the literature.

problem at hand, or use a generic strategy. As a common example of the latter, one may ensure that the initial centers are sufficiently far apart by picking $c_1$ at random, $c_2$ as far as possible from $c_1$, $c_3$ maximizing the sum of distances to $c_1$ and $c_2$ etc. One also typically runs the algorithm several times with random initial conditions and select the best solution over these multiple runs.

The second step of Algorithm 19.3 can be computationally challenging depending on the set $\mathcal{R}$ and the dissimilarity measure $\alpha$. When $\mathcal{R} = \mathbb{R}^d$ and $\alpha = \rho^2$ is the square Euclidean distance, the solution is explicit and $c_i$ is simply the average of all points in $A_i$. The resulting algorithm is the original incarnation of K-medoids, and called K-means [182, 121, 124]. K-means is probably the most popular clustering method and is often a step in more advanced approaches, as we will discuss later. The two steps of Algorithm 19.3 are then simplified as follows.

---

**Algorithm 19.4 (K-means)**
Let $T \subset \mathbb{R}^d$ be the training set. Start with an initial choice of $c_1, \ldots, c_p \in \mathbb{R}^d$ and iterate over the following two steps until stabilization:

(1) For $i = 1, \ldots, p$, let $A_i$ contain points $x \in T$ such that $|x - c_i|^2 = \min\{|x - c_j|^2, j = 1, \ldots, p\}$. In case of tie in this minimum, assign $x$ to only one of the tied sets (e.g., at random) to ensure that $A_1, \ldots, A_p$ is a partition.

(2) For $i = 1, \ldots, p$, let

$$c_i = \frac{1}{|A_i|} \sum_{x \in A_i} x$$

   if $A_i$ is not empty, or $c_i$ be a random point in $T$ otherwise.

---

### 19.3.2   Mixtures of Gaussian and deterministic annealing

Mixtures of Gaussian (MoG) were discusssed in chapter 16 and in Algorithm 16.2. Recall that they model the observed data $X$ together with a latent class variable $Z \in \{1, \ldots, p\}$ with joint distribution

$$f(x, z; \theta) = (2\pi)^{-\frac{d}{2}} (\det \Sigma_z)^{-\frac{1}{2}} \alpha_z e^{-\frac{1}{2}(x - c_z)^T \Sigma_z^{-1}(x - c_z)}$$

where $\theta$ contains the weights, $\alpha_1, \ldots, \alpha_p$, the means, $c_1, \ldots, c_p$ and the covariance matrices $\Sigma_1, \ldots, \Sigma_p$ (we create, hopefully without risk of confusion, a short-lived conflict of notation between the weights and the dissimilarity function). The posterior class probabilities

$$f_Z(i|x; \theta) = \frac{(\det \Sigma_i)^{-\frac{1}{2}} \alpha_i e^{-\frac{1}{2}(x - c_i)^T \Sigma_i^{-1}(x - c_i)}}{\sum_{j=1}^p (\det \Sigma_j)^{-\frac{1}{2}} \alpha_j e^{-\frac{1}{2}(x - c_j)^T \Sigma_j^{-1}(x - c_j)}}, \quad i = 1, \ldots, p,$$

which are computed in step 3 of Algorithm 16.2 can be interpreted as a likelihood that observation $x$ belongs to group $i$. As a consequence, the mixture of Gaussian algorithm can also be seen as a clustering method, in which one assigns each $x \in T$ to cluster $i$ when $i = \mathrm{argmax}\{f_Z(j|x,\theta) : j = 1,\ldots,p\}$, making an arbitrary decision in case of a tie.

In the special case in which all variances are fixed and equal to $\sigma^2 \mathrm{Id}_{\mathbb{R}^d}$, and all prior class probabilities are equal to $1/p$ (see remark 16.3), the EM algorithm for mixtures of Gaussian is also called "soft K-means", because it replaces the "hard" cluster assignments in K-means by "soft" ones represented by the update of the posterior distribution. We repeat its definition here for completeness (where $\theta = (c_1,\ldots,c_p)$).

---

**Algorithm 19.5 (Soft K-means)**
   1. Choose a number $\sigma^2 > 0$, a small constant $\epsilon$ and a maximal number of iterations $M$. Initialize the centers $c = (c_1,\ldots,c_p)$.

   2. At step $n$ of the algorithm, let $c$ be the current centers.

   3. Compute, for $x \in T$ and $i = 1,\ldots,p$

$$f_Z(i|x,\theta) = \frac{e^{-\frac{1}{2\sigma^2}|x-c_i|^2}}{\sum_{j=1}^{p} e^{-\frac{1}{2\sigma^2}|x-c_j|^2}}$$

and let $\zeta_i = \sum_{k=1}^{N} f_Z(i|x,\theta)$, $i = 1,\ldots,p$.

   4. For $i = 1,\ldots,p$, let

$$c_i' = \frac{1}{\zeta_i} \sum_{x \in T} x f_Z(i|x,\theta).$$

   5. If $|c' - c| < \epsilon$ or $n = M$: stop the algorithm.

   6. Replace $c$ by $c'$ and $n$ by $n+1$ and return to step 2.

---

When $\sigma^2 \to 0$, $f_Z(\cdot|x_k,\theta)$ converges to the uniform probability on indexes $j$ such that $c_j$ is closest to $x_k$, which is a Dirac measure unless there are ties. Class allocation and center updating become then asymptotically identical to the K-means algorithm. A variant of soft K-means, called deterministic annealing [169], applies Algorithm 19.5 while letting $\sigma$ slowly tend to 0. This new algorithm is experimentally more robust than K-means, in that it is less likely to be trapped in bad local minimums.

**Remark 19.4** The soft K-means algorithm can also be defined directly as an alternate minimization method for the objective function

$$F(c,f_Z) = \frac{1}{2} \sum_{x \in T} \sum_{j=1}^{p} f_Z(j|x)|x - c_j|^2 + \sigma^2 \sum_{x \in T} \sum_{j=1}^{p} f_Z(j|x) \log f_Z(j|x),$$

with the constraints $f_Z(j|x) \geq 0$ for all $j$ and $x$ and $\sum_{j=1}^{p} f_Z(j|x) = 1$. One can check (we leave this as an exercise) that Step 3 in Algorithm 19.5 provides the optimal $f_Z$ for $F$ when $c$ is fixed, and that Step 4 gives the optimal $c$ when $f_Z$ is fixed (see **??**). ◆

**Remark 19.5** We note that, if a K-means, soft K-means or MoG algorithm has been trained on a training set $T$, it is then easy to assign a new sample $\tilde{x}$ to one of the clusters. Indeed, for K-means, it suffices to determine the center closest to $\tilde{x}$, and for the other methods to maximize $f_Z(j|\tilde{x}, \theta)$, which is computable given the model parameters. In contrast, there was no direct way to do so using hierarchical clustering. ◆

### 19.3.3   Kernel (soft) K-means

We now consider the soft K-means algorithm in feature space, and introduce features $h_k = h(x_k)$ in an inner product space $H$ such that $\langle h_k, h_l \rangle_H = K(x_k, x_l)$ for some positive definite kernel. As usual, the underlying assumption is that the computation of $h(x)$ does not need to be feasible, while evaluations of $K(x, y)$ are easy. Let us consider the minimization of

$$\frac{1}{2} \sum_{x \in T} \sum_{j=1}^{p} f_Z(j|x) \|h(x) - c_j\|_H^2 + \sigma^2 \sum_{x \in T} \sum_{j=1}^{p} f_Z(j|x) \log f_Z(j|x)$$

for some $\sigma^2 > 0$ (kernel K-means corresponds to taking the limit $\sigma^2 \to 0$). Given $f_Z$, the optimal centers are

$$c_j = \frac{1}{\zeta_j} \sum_{x \in T} f_Z(j|x) h(x)$$

with $\zeta = \sum_{x \in T} f_Z(j|x)$. They belong to the feature space, $H$, and are therefore not computable in general. However, the distance between them and a point $h(y) \in H$ is explicit and given by

$$\|h(y) - c_j\|_H^2 = K(y, y) - \frac{2}{\zeta_j} \sum_{x \in T} f_Z(j|x) K(y, x) + \frac{1}{\zeta_j^2} \sum_{x, x' \in T} f_Z(j|x) f_Z(j|x') K(x, x').$$

The class probabilities at each iteration can therefore be updated using

$$f_Z(j|x) = \frac{e^{-\|h(x) - c_j\|_H^2 / 2\sigma^2}}{\sum_{j'=1}^{p} e^{-\|h(y) - c_{j'}\|_H^2 / 2\sigma^2}}.$$

This yields the soft kernel K-means algorithm, that we repeat below.

**Algorithm 19.6 (Kernel soft K-means)**
Let $T \subset \mathbb{R}^d$ be the training set. Initialize the algorithm with some choice for $f_Z(j|x)$, $j = 1, \ldots, p$, $x \in T$ (for example: $f_Z(j|x) = 1/p$ for all $j$ and $x$).

(1) For $j = 1, \ldots, p$ and $x \in T$ compute

$$\|h(x) - c_j\|_H^2 = K(x,x) - \frac{2}{\zeta_j} \sum_{x' \in T} f_Z(j|x') K(x,x') + \frac{1}{\zeta_j^2} \sum_{x',x'' \in T} f_Z(j|x') f_Z(j|x'') K(x',x'')$$

with $\zeta_j = \sum_{x' \in T} f_Z(j|x')$.

(2) Compute, for $x \in T$ and $j = 1, \ldots, p$,

$$f_Z(j|x) = \frac{e^{-\|h(x) - c_j\|_H^2 / 2\sigma^2}}{\sum_{j'=1}^{p} e^{-\|h(y) - c_{j'}\|_H^2 / 2\sigma^2}} .$$

(3) If the variation of $f_Z$ compared to the previous iteration is small, or if a maximum number of iterations has been reached, exit the algorithm.

(4) Return to step 1.

After convergence, the clusters are computed by assigning $x$ to $A_i$ when $i = \mathrm{argmax}\{f_Z(j|x) : j = 1, \ldots, p\}$, making an arbitrary decision in case of a tie.

---

For "hard" K-means (with $\sigma^2 \to 0$), step 2 simply updates $f_Z(j|x)$ as the uniform probability on the set of indexes $j$ at which $\|h(x) - c_j\|_H^2$ is minimal.

### 19.3.4 Convex relaxation

We return to the initial formulation of K-means for Euclidean data, as a minimization, over all partitions $\mathcal{A} = \{A_1, \ldots, A_K\}$ of $\{1, \ldots, N\}$ of

$$W(\mathcal{A}) = \sum_{j=1}^{K} \sum_{k \in A_j} |x_k - c_j|^2$$

where $c_j$ is the average of the points $x_j$ such that $j \in A_j$. We start with a simple transformation expressing this function in terms of the matrix $S_\alpha$ of square distances

$\alpha(x_k, x_l) = |x_k - x_l|^2$. Indeed, we have

$$\sum_{k \in A_j} |x_k - c_j|^2 = \sum_{k \in A_j} |x_k|^2 - \frac{1}{|A|} \left| \sum_{k \in A_j} x_k \right|$$

$$= \sum_{k \in A_j} |x_k|^2 - \frac{1}{|A|} \sum_{k,l \in A_j} x_k^T x_l$$

$$= \frac{1}{2|A_j|} \sum_{k,l \in A_j} (|x_k|^2 + |x_l|^2 - 2x_k^T x_l)$$

$$= \frac{1}{2|A_j|} \sum_{k,l \in A_j} |x_k - x_l|^2$$

Introduce the vector $u_j \in \mathbb{R}^N$ with coordinates $u_j^{(k)} = 1/\sqrt{|A_j|}$ for $k \in A_j$ and 0 otherwise. Then

$$\frac{1}{2|A_j|} \sum_{k,l \in A_j} |x_k - x_l|^2 = \frac{1}{2} u_j^T S_\alpha u_J = \frac{1}{2} \text{trace}(S_\alpha u_j u_j^T). \tag{19.4}$$

Let

$$Z(\mathcal{A}) = \sum_{j=1}^{p} u_j u_j^T,$$

so that $Z(\mathcal{A})$ has entries $Z^{(k,l)}(\mathcal{A}) = 1/|A_j|$ for $k, l \in A_j$, $j = 1, \ldots p$ and 0 for all other $k, l$. Summing (19.4) over $j$, we get

$$W(\mathcal{A}) = \frac{1}{2} \text{trace}(S_\alpha Z(\mathcal{A})).$$

The matrix $Z(\mathcal{A})$ is symmetric, has non-negative entries. It moreover satisfies $Z(\mathcal{A})\mathbb{1}_N = \mathbb{1}_N$ and $Z(\mathcal{A})^2 = Z(\mathcal{A})$. Interestingly, these properties characterize matrices $Z$ associated with partitions, as stated in the next proposition [153, 152].

**Proposition 19.6** *Let $Z \in \mathcal{M}_N(\mathbb{R})$ be a symmetric matrix with non-negative entries satisfying $Z\mathbb{1}_N = \mathbb{1}_N$ and $Z^2 = Z$. The there exists a partition $\mathcal{A}$ of $\{1, \ldots, N\}$ such that $Z = Z(\mathcal{A})$.*

Proof  Note that $Z$ being symmetric and satisfying $Z^2 = Z$ imply that it is an orthogonal projection with eigenvalues 0 and 1. In particular $Z$ is positive semidefinite. This implies that, for all $i, j \in \{1, \ldots, N\}$, one has

$$Z(i,j)^2 \leq Z(i,i), Z(j,j).$$

This inequality combined with $\sum_{j=1}^{N} Z(k, j) = 1$ (expressing $Z\mathbb{1}_N = \mathbb{1}_N$) shows that all diagonal entries of $Z$ are positive.

Define on $\{1,\dots,N\}$ the relation $k \sim j$ if and only if $Z(j,k) > 0$. The relation is symmetric and we just checked that $k \sim k$ for all $k$. It is also transitive, from the relation (deriving from $Z^2 = Z$)

$$Z(k,j) = \sum_{i=1}^{N} Z(k,i)Z(i,j)$$

which shows (since all terms in the sum are non-negative) that $k \sim i$ and $j \sim i$ imply $k \sim j$.

Let $\mathcal{A} = \{A_1,\dots,A_q\}$ be the partition of $\{1,\dots,N\}$ formed by the equivalence classes for this relation. We now show that $Z = Z(\mathcal{A})$.

We have, for all $k,j \in \{1,\dots,N\}$

$$\sum_{i=1}^{N} Z(k,i)(Z(k,j) - Z(i,j)) = Z(k,j)\sum_{i=1}^{N} Z(k,i) - \sum_{i=1}^{N} Z(k,i)Z(i,j)$$

$$= Z(k,j) - \sum_{i=1}^{N} Z(k,i)Z(i,j) = 0$$

Now, if $k,j \in A_s$ for some $s$, the identity reduces to

$$\sum_{i \in A_s} Z(k,i)(Z(k,j) - Z(i,j)) = 0. \tag{19.5}$$

Choose $k$ such that $Z(k,k) = \max\{Z(i,i) : i \in A_s\}$. Then, for all $i,j \in A_s$, $Z(i,j) \leq \sqrt{Z(i,i)Z(j,j)} \leq Z(k,k)$ and (19.5) for $j = k$ yields

$$\sum_{i \in A_s} Z(k,i)(Z(k,k) - Z(k,i)) = 0,$$

which is only possible (since all $Z(k,i)$ are positive) if $Z(k,i) = Z(k,k)$ for all $i \in A_s$. From $Z(k,i) \leq \sqrt{Z(i,i)Z(k,k)}$, we get $Z(i,i) = Z(k,k)$ for all $i$, and therefore (reapplying what we just found to $i$ instead of $k$) $Z(i,j) = Z(i,i) = Z(k,k)$ for all $i,j \in A_s$. Finally, we have

$$1 = \sum_{i \in A_s} Z(k,i) = |A_s|Z(k,k)$$

showing that $Z(k,k) = 1/|A_s|$ and completing the proof that $Z = Z(\mathcal{A})$. ∎

Note that the number of clusters, $|\mathcal{A}|$ is equal to the trace of $Z(\mathcal{A})$. This shows that minimizing $W(\mathcal{A})$ over partitions with $p$ clusters is equivalent to the constrained optimization problem minimizing

$$G(Z) = \text{trace}(S_\alpha Z) \tag{19.6}$$

over all matrices $Z$ such that $Z \geq 0$, $Z^T = Z$, $Z\mathbb{1}_N = \mathbb{1}_N$, trace$(Z) = p$ and $Z^2 = Z$. This is still a difficult problem, since it is equivalent to K-means, which is NP hard. Seeing the problem in this form, however, is more amenable to approximations and, in particular, convex relaxations.

In [152], it is proposed to use a semidefinite program (SDP) as a relaxation. The conditions $Z = Z^T$ and $Z^2 = Z$ require that all eigenvalues of $Z$ are either 0 or 1, and a direct relaxation is to replace these constraints by $Z^T = Z$ and $0 \leq Z \leq \mathrm{Id}_{\mathbb{R}^N}$. The last inequality is however redundant if we add the conditions [2] $Z \geq 0$ and $Z\mathbb{1} = \mathbb{1}$. This is a consequence of the Perron-Frobenius theorem which states that a matrix $\tilde{Z}$ with positive entries has a largest (in modulus) real eigenvalue, which has multiplicity one and is associated with an eigenvector with positive coordinates, the latter eigenvector being (up to multiplication by a constant) the unique eigenvector of $\tilde{Z}$ with positive coordinates. So, if a matrix $\tilde{Z}$ is symmetric, satisfies $\tilde{Z} > 0$ and $\tilde{Z}\mathbb{1}_N = \mathbb{1}_N$, then $\tilde{Z} \leq \mathrm{Id}_{\mathbb{R}^N}$. Applying this result to $\tilde{Z} = (1 - \epsilon)Z + (\epsilon/N)\mathbb{1}_N\mathbb{1}_N^T$ and letting $\epsilon$ tend to 0 shows that any matrix $Z$ with non-negative entries satisfying $Z\mathbb{1}_N = \mathbb{1}_N$ also satisfies $Z \leq \mathrm{Id}_{\mathbb{R}^N}$.

This provides the following SDP relaxation of K-means [152]: minimize

$$G(Z) = \mathrm{trace}(S_\alpha Z) \tag{19.7}$$

subject to $Z^T = Z$, $Z\mathbb{1}_N = \mathbb{1}_N$, trace$(Z) = p$, $Z \geq 0$, $Z \succeq 0$.

Clusters can be immediately inferred from the columns of the matrix $Z(\mathcal{A})$, since they are identical for two indices in the same cluster, and orthogonal to each other for two indices in different clusters. Let $z_1(\mathcal{A}), \ldots, z_N(\mathcal{A})$ denote the columns of $Z(\mathcal{A})$ and $\bar{z}_k(\mathcal{A}) = z_k(\mathcal{A})/|z_k(\mathcal{A})|$. One has $|\bar{z}_k(\mathcal{A}) - \bar{z}_l(\mathcal{A})| = 0$ if $k$ and $l$ belong to the same cluster and $\sqrt{2}$ otherwise.

These properties will not necessarily be satisfied by a solution, say, $Z^*$, of the SDP relaxation, but, assuming that the approximation is good enough, one may still consider the normalized columns of $Z^*$ and expect them to be similar for indices in the same cluster, and away from each other otherwise. Denoting by $\bar{z}_1^*, \ldots, \bar{z}_N^*$ these normalized columns, one can then run on them the standard K-means algorithm, or a spectral clustering method such as those described in the next sections, to infer clusters.

**Remark 19.7** Clearly, one can use any symmetric matrix $S$ in the definition of $G$ in (19.6) and (19.7). The method is equivalent to, or to a relaxation of, K-means only when $S$ is formed with squared norms in inner-product spaces, which does include kernel K-means, for which

$$\alpha(x_k, x_l) = K(x_k, x_k) - 2K(x_k, x_l) + K(x_l, x_l).$$

---

[2]Recall that $Z \succeq 0$ means that $Z$ is positive definite, while $Z \geq 0$ indicates that all its entries are non-negative.

If $\alpha$ is an arbitrary discrepancy measure, the minimization of $G(Z)$ still makes sense, since it is equivalent to minimizing

$$G(Z(\mathcal{A})) = \sum_{j=1}^{p} D_\alpha(A_j).$$

where

$$D_\alpha(A) = \frac{1}{|A|} \sum_{x,y \in A} \alpha(x,y). \tag{19.8}$$

is a (normalized) measure of size, that we will call the $\alpha$-*dispersion* of a finite set $A$.♦

**Remark 19.8** Instead of using dissimilarities, some algorithms are more naturally defined in terms of similarities. Given such a similarity measure, say, $\beta$, one must maximize rather than minimize the index $\Delta_\beta$ (which becomes, rather than a measure of dispersion, a measure of concentration).

One passes from a dissimilarity $\alpha$ to a similarity $\beta$ by applying a decreasing function to the former, a common choice being

$$\beta(x,x') = \exp(-\alpha(x,x')/\tau)$$

for some $\tau > 0$.

Alternatively, one can fix an element $x_0 \in \mathcal{R}$ and let

$$\beta(x,y) = \alpha(x,x_0) + \alpha(y,x_0) - \alpha(x,y) - \alpha(x_0,x_0),$$

(note that the last term, $\alpha(x_0,x_0)$ is generally equal to 0). For example, if $\alpha(x,y) = |x-y|^2$, then $\beta(x,y) = 2(x-x_0)^T(y-x_0)$ (for which it is natural to take $x_0 = 0$). If $\alpha$ is a distance (not squared!), then $\beta \geq 0$ by the triangular inequality. In this case, we have

$$\Delta_\beta(A_1,\ldots,A_p) = \sum_{k=1}^{n} D_\beta(A_k)$$

$$= \sum_{k=1}^{p} \frac{1}{|A_p|} \sum_{x,y \in A_k} \alpha(x,x_0) + \sum_{k=1}^{p} \frac{1}{|A_p|} \sum_{x,y \in A_k} \alpha(y,x_0)$$

$$- \sum_{k=1}^{p} \frac{1}{|A_p|} \sum_{x,y \in A_k} \alpha(x_0,x_0) - \sum_{k=1}^{p} \frac{1}{|A_p|} \sum_{x,y \in A_k} \alpha(x,x_0)$$

$$= 2 \sum_{k=1}^{p} \sum_{x \in A_k} \alpha(x,x_0) - \sum_{k=1}^{p} |A_k|\alpha(x_0,x_0) - \Delta_\alpha(A_1,\ldots,A_p)$$

$$= 2 \sum_{x \in T} \alpha(x,x_0) - |T|\alpha(x_0,x_0) - \Delta_\alpha(A_1,\ldots,A_p) \qquad ♦$$

so that minimizing $\Delta_\alpha$ is equivalent to maximizing $\Delta_\beta$.

## 19.4   Spectral clustering

### 19.4.1   Spectral approximation of minimum discrepancy

One refers to spectral methods algorithms that rely on computing eigenvectors and eigenvalues (the spectrum) of data-dependent matrices. In the case of minimizing discrepancies, they can be obtained by further simplifying (19.7), essentially by removing constraints.

One indeed gets a simpler problem if the non-negativity constraint, $Z \geq 0$, is removed. Doing so, one cannot guarantee anymore that $Z \preceq \mathrm{Id}_{\mathbb{R}^N}$, so we need to reinstate this constraint. We will first make the further simplification to remove the constraint $Z \mathbb{1}_N = \mathbb{1}_N$, the problem becoming minimizing $\mathrm{trace}(S_\alpha Z)$ over all $Z \in \mathcal{S}_N^+(\mathbb{R})$ such that $0 \preceq Z \preceq \mathrm{Id}_{\mathbb{R}^N}$ and $\mathrm{trace}(Z) = p$. Decomposing $Z$ in an eigenbasis, i.e., looking for it in the form

$$Z = \sum_{j=1}^{N} \xi_j e_j e_j^T,$$

this is equivalent to minimizing

$$\sum_{j=1}^{N} \xi_j e_j^T S_\alpha e_j \tag{19.9}$$

subject to $0 \leq \xi_j \leq 1$, $\sum_{j=1}^{N} \xi_j = p$ and $u_1, \ldots, u_N$ orthonormal basis of $\mathbb{R}^N$. First consider minimization with respect to the basis, fixing $\xi$. There is obviously no loss of generality in requiring that $\xi_1 \leq \xi_2 \leq \cdots \leq \xi_N$, and using corollary 2.4 (adapted to minimizing (19.9) rather than maximizing it) we know that an optimal basis is given by the eigenvectors of $S_\alpha$, ordered with non-decreasing eigenvalues. Letting $\lambda_1 \leq \cdots \leq \lambda_N$ denote these eigenvalues, we find that $\xi_1, \ldots, x_N$ must be a non-decreasing sequence minimizing

$$\sum_{j=1}^{N} \lambda_j \xi_j$$

subject to $0 \leq \xi_k \leq 1$ and $\sum_{j=1}^{N} \xi_j = p$. The optimal solution is obtained by taking

$\xi_1 = \cdots = \xi_p = 1$, since, for any other solution

$$\sum_{j=1}^{N} \lambda_j \xi_j - \sum_{j=1}^{p} \lambda_j \geq \lambda_{p+1} \sum_{j=p+1}^{N} \xi_j + \sum_{j=1}^{p} \lambda_j (\xi_j - 1)$$

$$= \lambda_{p+1} \sum_{j=1}^{p} (1 - \xi_j) + \sum_{j=1}^{p} \lambda_j (\xi_k - 1)$$

$$= \sum_{j=1}^{p} (\lambda_{p+1} - \lambda_j)(1 - \xi_j)$$

$$\geq 0.$$

The following algorithm (similar to that discussed in [64]) summarizes this discussion.

---

**Algorithm 19.7 (Spectral clustering: version 1)**
Let $S_\alpha$ be an $N \times N$ discrepancy matrix. Let $p$ denote the number of clusters.

(1) Compute the eigenvectors of $S_\alpha$ associated with the $p$ smallest eigenvalues.

(2) Denoting these eigenvectors by $e_1, \ldots, e_p$, define $y_1, \ldots, y_N \in \mathbb{R}^p$ by $y_k^{(j)} = e_j^{(k)}$.

(3) Run K-means on $(y_1, \ldots, y_N)$ to determine a partition.

---

This algorithm needs to be slightly modified if one also wants $Z$ to satisfy $Z\mathbb{1} = \mathbb{1}$. In that case, $\mathbb{1}$ is one of the eigenvectors (with eigenvalue 1), and the others are orthogonal to it. As a consequence, one now looks for $Z$ in the form

$$Z = \sum_{k=1}^{N-1} \xi_j e_j e_j^T + \frac{1}{N} \mathbb{1}\mathbb{1}^T$$

leading to the minimization of

$$\sum_{j=1}^{N-1} \xi_j e_j^T S_\alpha e_j + \frac{1}{N} \mathbb{1}^T S_\alpha \mathbb{1}$$

over all $\xi_1, \ldots, \xi_{N-1}$ such that $0 \leq \xi_j \leq 1$ and $\sum_{j=1}^{N} \xi_j = p - 1$, and over all $e_1, \ldots, e_{N-1}$ such that $e_1, \ldots, e_{N-1}, \mathbb{1}/\sqrt{N}$ form an orthonormal basis. The main difference with the previous problem is that we now need to ensure that all $e_j$ are perpendicular to $\mathbb{1}$.

To achieve this, introduce the projection matrix $P = \mathrm{Id}_{\mathbb{R}^N} - \mathbb{1}\mathbb{1}^T/N$ and let $\tilde{S}_\alpha = PS_\alpha P$. Then, since $u^T\mathbb{1} = 0$ implies $u^T\tilde{S}_\alpha u = u^T S_\alpha u$, it is equivalent to minimize

$$\sum_{j=1}^{N-1} \xi_j e_j^T \tilde{S}_\alpha e_j$$

over all $\xi_1,\ldots,\xi_{N-1}$ such that $0 \leq \xi_j \leq 1$ and $\sum_{j=1}^{N} \xi_j = p - 1$, and over all $e_1,\ldots,e_{N-1}$ such that $e_1,\ldots,e_{N-1}, \mathbb{1}/\sqrt{N}$ form an orthonormal basis. Because $\tilde{S}_\alpha \mathbb{1} = 0$, we know that $\tilde{S}_\alpha$ can be diagonalized in an orthonormal basis $(e_1,\ldots,e_{N-1}, \mathbb{1}/\sqrt{N})$, and we obtain an optimal solution by selecting the $p - 1$ vectors associated with smallest eigenvalues, with associated $\xi_j = 1$. We therefore get a modified version of the spectral clustering algorithm.

---

**Algorithm 19.8 (Spectral clustering: version 2)**
Let $S_\alpha$ be an $N \times N$ discrepancy matrix. Let $p$ denote the number of clusters. Let $P = \mathrm{Id}_{\mathbb{R}^N} - \mathbb{1}_N \mathbb{1}_N^T/N$.

(1) Compute $\tilde{S}_\alpha = PS_\alpha P$

(2) Compute the eigenvectors of $\tilde{S}_\alpha$ associated with the $p - 1$ smallest eigenvalues.

(3) Denoting these eigenvectors by $e_1,\ldots,e_{p-1}$, define $y_1,\ldots,y_N \in \mathbb{R}^{p-1}$ by $y_k^{(j)} = e_j^{(k)}$.

(4) Run K-means on $(y_1,\ldots,y_N)$ to determine a partition.

---

## 19.5   Graph partitioning

Similarity measures are often associated with graph structures, with a goal of finding a partition of their set of vertices. So, let $T$ denote the set of these vertices and assume that to all pairs $x, y \in T$, one attribute a weight given by $\beta(x,y)$, where $\beta$ is assumed to be non-negative. We define $\beta$ for all $x, y \in T$, but we interpret $\beta(x,y) = 0$ as marking the absence of an edge between $x$ and $y$.

Let $V$ denote the vector space of all functions $f : T \to \mathbb{R}$ (we have $\dim(V) = |T|$). This space can be equipped with the standard Euclidean norm, that we will call in this section the $L^2$ norm (by analogy with general spaces of square integrable functions), letting,

$$|f|_2^2 = \sum_{x \in T} f(x)^2.$$

One can also associate a measure of smoothness for a function $f \in V$ by computing the discrete "$H^1$" semi-norm,

$$|f|_{H^1}^2 = \sum_{x,y \in T} \beta(x,y)(f(x) - f(y))^2.$$

With this definition, "smooth functions" tend to have similar values at points $x, y$ in $T$ such that $\beta(x, y)$ is large while there is less constraint when $\beta(x, y)$ is small. In particular, $|f|_{H^1} = 0$ if and only if $f$ is constant on connected components of the graph.[3]

The notion of connected components, combined with thresholding, can be used to build a hierarchical family of partitions of the graph. Define, for all $t > 0$, the thresholded weights $\beta^{(t)}(x, y) = \max(\beta(x, y) - t, 0)$. The set of connected components associated with the pair $(V, \beta^{(t)})$ forms a partition, say, $\mathcal{A}^{(t)}$, of $T$. The resulting set of partitions is *nested* in the sense that, if $s < t$, the sets forming the partition $\mathcal{A}^{(s)}$ are unions of sets forming $\mathcal{A}^{(t)}$. This thresholding procedure is not always satisfactory, however, because there does not always exist a fixed value of $t$ that produces a good quality cluster decomposition.

If there exists $p$ connected components, then the subspace of all functions $f \in V$ such that $|f|_{H_1} = 0$ has dimension $p$. If $C_1, \ldots, C_p$ are the connected components, this space is generated by the functions $\delta_{C_k}$, $k = 1, \ldots, p$, with $\delta_{C_k}(x) = 1$ if $x \in C_k$ and $0$ otherwise. These functions form, in addition, an orthogonal system for the Euclidean inner product: $\langle \delta_{C_k}, \delta_{C_l} \rangle_2 = 0$ if $k \neq l$.

One can write $\frac{1}{2}|f|^2_{H^1} = f^T L f$ where $L$, called the *Laplacian operator* associated to the considered graph, is defined by

$$Lf(x) = \sum_{y \in T} L(x, y) f(y)$$

and

$$L(x, y) = \left( \sum_{z \in T} \beta(x, z) \right) \mathbf{1}_{x=y} - \beta(x, y). \tag{19.10}$$

The vectors $\delta_{C_k}$, $k = 1, \ldots, p$ are then an orthogonal basis of the null space of $L$. Conversely, let $(e_1, \ldots, e_p)$ be any basis of this null space. Then, there exists an invertible matrix $A = (a_{ij}, i, j = 1, \ldots, p)$ such that

$$e_i(x) = \sum_{j=1}^{p} a_{ij} \delta_{C_j}(x).$$

Associate to each $x \in T$ the vector $e(x) = \begin{pmatrix} e_1(x) \\ \vdots \\ e_p(x) \end{pmatrix} \in \mathbb{R}^p$. Then, for any $x, y \in T$, we have

$e(x) = e(y)$ if and only if $\delta_{C_j}(x) = \delta_{C_j}(y)$ for all $j = 1, \ldots, p$ (because $A$ is invertible),

---

[3]Two nodes $x$ and $y$ are connected in the graph if there is a sequence $z_0, \ldots, z_n$ in $T$ such that $z_0 = x$, $z_n = y$ and $\beta(z_i, z_{i-1}) > 0$ for $i = 1, \ldots, n$. This provides an equivalence relation and equivalent classes are called connected components.

that it, if and only if $x$ and $y$ belong to the same connected component. So, given any basis of the null space of $L$, the function $x \mapsto e(x)$ determines these connected components. So, a—not very efficient—way of determining the connected components of the graph can be to diagonalize the operator $L$ (written as an $N$ by $N$ matrix, where $N = |T|$), extract the $p$ eigenvectors $e_1,\ldots,e_p$ associated with eigenvalue zero and deduce from the function $e(x)$ above the set of connected components.

Now, in practice, the graph associated to $T$ and $\beta$ will not separate nicely into connected components in order to cluster the training set. Most of the time, because of noise or some weak connections, there will be only one such component, or in any case much less than what one would expect when clustering the data. The previous discussion suggests, however, that in the presence of moderate noise in the connection weights, one may expect that the eigenvectors associated to the $p$ smallest eigenvalues of $L$ provide vectors $e(x), x \in T$ such that $e(x)$ and $e(y)$ have similar values if $x$ and $y$ belong to the same cluster (see 19.2). In such cases, these clusters should be easy to determine using, say, K-means on the transformed dataset $\tilde{T} = (e(x), x \in T)$. This is summarized in the following algorithm.

---

**Algorithm 19.9 (Spectral Graph Partitioning)**
Let $T \subset \mathcal{R}$ be the training set and $(x,y) \mapsto \beta(x,y)$ a similarity measure defined on $T \times T$. Let $p$ be the desired number of clusters.

(1) Form the Laplacian operator described in (19.10) and let $e_1,\ldots,e_p$ be its eigenvectors associated to the $p$ lowest eigenvalues. For $x \in T$, let $e(x) \in \mathbb{R}^p$ be given by

$$e(x) = (e_1(x),\ldots,e_p(x))^T \in \mathbb{R}^p.$$

(2) Apply the K-means algorithm (or one of its variants) with $p$ clusters to $\tilde{T} = (e(x), x \in T)$.

---

## 19.6   Deciding the number of clusters

### 19.6.1   Detecting elbows

The number, $p$, of subsets with respect to which the population should be partitioned is rarely known a priori, and several methods have been introduced in the literature in order to assess the ideal number of clusters. We now review some of these methods, and denote, for this purpose, by $\mathcal{L}^*(p)$ the minimized cost function obtained with $p$ clusters, e.g., using (19.3),

$$\mathcal{L}^*(p) = \min\{\mathbb{W}_\alpha(A_1,\ldots,A_p,c_1,\ldots,c_p) : A_1,\ldots,A_p \text{ partition of } T, c_1,\ldots,c_p \in \mathcal{R}\},$$

Figure 19.2: Example of data transformed using the eigenvectors of the graph Laplacian. Left: Original data. Center: Result of a Kmeans algorithm with three clusters applied to the transformed data (2D projection). Right: Visualization of the cluster labels on the original data.

in the case of K-medoids (this definition is algorithm dependent). It is clear that $\mathcal{L}^*$ is a decreasing function of $p$. It is also natural to expect that $\mathcal{L}^*$ should decrease significantly when $p$ is smaller than the correct number of clusters, while the variation should be more marginal when $p$ is overestimated, because the cost in putting together two sets of points that are far apart (which happens when $p$ is too small) is typically larger than the gain in splitting a homogeneous region in two.

The simplest approach in this context is to visualize $\mathcal{L}^*(p)$ as a function of $p$ and try to locate at which value the resulting curve makes an "elbow," i.e., switches from a sharply decreasing slope to a milder one. Figure 19.3 provides an illustration of this visualization when the true number of clusters is three (the data in each cluster following a normal distribution). When the clusters are well separated, an elbow clearly appears on the graph of $\Gamma_\alpha^*$, but this situation is harder to observe when clusters overlap with each other.

One can measure the "curvature" at the elbow using the distance between each point in the graph of $(p, W_\alpha^*(p))$ and the line between its predecessor and successor. The result gives the criterion

$$C(p) = \frac{\mathcal{L}^*(p+1) + \mathcal{L}^*(p-1) - 2\mathcal{L}^*(p)}{\sqrt{(\mathcal{L}^*(p+1) - \mathcal{L}^*(p-1))^2 + 4}},$$

specifying the elbow point as the value of $p$ at which $C$ attains its maximum. For both examples in fig. 19.3, this method returns the correct number of clusters (3).

### 19.6.2 The Caliński and Harabasz index

Several other criteria have been introduced in the literature. Caliński and Harabasz [46] propose to minimize the ratio of normalized between-group and within-groups sums of squares associated with K-means. For a given $p$, let $c_1, \ldots, c_p$ denote the optimal centers, and $A_1, \ldots, A_p$ the optimal partition, with $N_k = |A_k|$. The normalized

Figure 19.3: Elbow graphs for K-means clustering for two populations generated as mixtures of Gaussian.

between-group sum of squares is

$$h_\alpha(p) = \frac{1}{p-1} \sum_{k=1}^{p} N_k |c_k - \overline{x}|^2$$

and the normalized within-group sum of squares is

$$w_\alpha(p) = \frac{1}{N-p} \sum_{k=1}^{p} \sum_{x \in A_k} |x - c_k|^2$$

Caliński and Harabasz [46] suggest to maximize $\gamma_{CH}(p) = h_\alpha(p)/w_\alpha(p)$.

This criterion can be extended to other types of cluster analysis. We have seen in section 19.4 that, when $\alpha(x, y) = |x - y|^2$,

$$\frac{1}{2} \sum_{k=1}^{p} \sum_{x,y \in A_k} \alpha(x, y)/N_k = \sum_{k=1}^{p} \sum_{x \in A_k} |x - c_k|^2.$$

We also have

$$\sum_{x \in T} |x - \overline{x}|^2 = \sum_{k=1}^{p} \sum_{x \in A_k} |x - c_k|^2 + \sum_{k=1}^{p} N_k |c_k - \overline{x}|^2$$

and the left-hand side is also equal to

$$\frac{1}{2N} \sum_{x,y \in T} \alpha(x,y).$$

It follows that, when $\alpha(x,y) = |x-y|^2$,

$$h_\alpha(p) = \frac{1}{2(p-1)} \left( \frac{1}{N} \sum_{x,y \in T} \alpha(x,y) - \sum_{k=1}^{p} \sum_{x,y \in A_k} \alpha(x,y)/N_k \right)$$

and

$$w_\alpha(p) = \frac{1}{2(N-p)} \sum_{k=1}^{p} \sum_{x,y \in A_k} \alpha(x,y)/N_k.$$

These expressions can obviously be applied to any dissimilarity measure, extending $\gamma_{CH}$ to general clustering problems.

### 19.6.3 The "silhouette" index

For $x \in T$, let

$$d_\alpha(x, A_k) = \frac{1}{N_k} \sum_{y \in A_k} \alpha(x,y).$$

Let $a_\alpha(x,p) = d_\alpha(x, A(x))$ and $b(x,p) = \min\{d_\alpha(x, A_k) : A_k \neq A(x)\}$. Define the silhouette index of $x$ in the segmentation [170]by

$$s_\alpha(x,p) = \frac{b_\alpha(x,p) - a_\alpha(x,p)}{\max(b_\alpha(x,p), a_\alpha(x,p))} \in [-1, 1].$$

This index measures how well $x$ is classified in the partitioning. It is large when the mean distance between $x$ and other objects in its class is small compared to the minimum mean distance between $x$ and any other class. In order to estimate the best number of clusters with this criterion, one then can maximize the average index:

$$\gamma_R(p) = \frac{1}{N} \sum_{x \in T} s_\alpha(x,p).$$

**Remark 19.9** One can rewrite the Caliński and Harabasz index using the notation introduced for the silhouette index. Indeed, let $A(x)$ be the cluster $A_k$ to which $x$ belongs. Then

$$h_\alpha(p) = \frac{1}{2(p-1)} \sum_{x \in T} \sum_{k=1}^{p} \frac{N_k}{N} (d_\alpha(x, A_k) - d_\alpha(x, A(x)))$$

and

$$w_\alpha(p) = \frac{1}{2(N-p)} \sum_{k=1}^{p} \sum_{x \in A_k} d_\alpha(x, A_k).$$

♦

Figure 19.4: Division of the unit square into clusters for uniformly distributed data.

### 19.6.4   Comparing to homogeneous data

Several selection methods choose $p$ based on the comparison of the data to a "null hypothesis" of no cluster. For example, assume that K-means is applied to a training set $T$ where samples are drawn uniformly according to the uniform distribution on $[0,1]^d$. Given centers, $c_1, \ldots, c_p$, let $\bar{A}_k$ be the set of points in $[0,1]^d$ that are closer to $c_k$ than to any other point. Then the segmentation of $T$ is formed by the sets $A_k = \{x \in T : x \in \bar{A}_k\}$ and, for large enough $N$, we can approximate $|A_k|/N$ (by the Law of Large Numbers) by the volume of the set $\bar{A}_k$, that we will denote by $\mathrm{vol}(\bar{A}_k)$.

Let us assume that $c_1, \ldots, c_p$ are uniformly spaced, so that the sets $\bar{A}_k$ have similar volumes (close to $1/p$) and have roughly spherical shapes (see fig. 19.4). This implies that

$$\int_{\bar{A}_k} |x - c_k|^2 dx \simeq \mathrm{vol}(A_k) \frac{r_p^2 d}{d+2}$$

where $r_p$ is the radius of a sphere of volume $1/p$, i.e., $pr_p^d \simeq d/\Gamma_{d-1}$ where $\Gamma_{d-1}$ is the surface area of the unit sphere in $\mathbb{R}^d$. So, we should have, for some constant $C$ that only depends on $d$,

$$\sum_{x \in A_k} |x - c_k|^2 \simeq N_k \int_{\bar{A}_k} |x - c_k|^2 dx \simeq C(d)(pN)p^{-2/d-1} = C(d)Np^{-2/d}.$$

This suggests that, for fixed $N$ and $d$, $p^{2/d}\mathcal{L}^*(p)$ should vary slowly when $p$ overesti-

mate the number of clusters (assuming that this operation divides an homogeneous cluster). Based on this analysis, Krzanowski and Lai [111] introduced the difference-ratio criterion, namely,

$$\gamma_{KL}(p) = \left| \frac{(p-1)^{\frac{2}{d}}\mathcal{L}^*(p-1) - p^{\frac{2}{d}}\mathcal{L}^*(p)}{p^{\frac{2}{d}}\mathcal{L}^*(p) - (p+1)^{\frac{2}{d}}\mathcal{L}^*(p+1)} \right|,$$

and estimate the number of clusters by taking $p$ maximizing $\gamma_{KL}$.

Another similar approach, introduced by Sugar and James [185], is based on an analysis of mixtures of Gaussian, namely assuming an underlying model with $p_0$ groups, where data in group $k$ follow a Gaussian distribution $\mathcal{N}(\mu_k, \text{Id})$ (possibly after standardizing the covariance matrix). In that work, the authors show that, if $d$ (the dimension) tends to infinity, with the minimal distance between centers growing proportionally to $\sqrt{d}$, then $\mathcal{L}^*(p)/d$ tends to infinity when $p < p_0$. They also show that, with similar assumptions, $\mathcal{L}^*(p)/d$ behaves like $p^{-2/d}$ for $p \geq p_0$, still for large dimensions. Based on this, they suggest using the criterion

$$\gamma_{SJ}(p) = \left( \frac{\mathcal{L}^*(p)}{d} \right)^{-\nu} - \left( \frac{\mathcal{L}^*(p-1)}{d} \right)^{-\nu}$$

(with the convention that $\mathcal{L}^*(0) = 0$) for some positive number $\nu$ and select the value of $p$ that maximizes $\gamma_{SJ}$. Indeed, in the case of Gaussian mixtures, the choice $\nu = d/2$ ensures that, in large dimensions, $\gamma_{SJ}(p)$ is small for $p < p_0$, that it is close to 1 for $p > p_0$ and close to $p_0$ for $p = p_0$.

A more computational approach, based on Monte-Carlo simulations has been introduced in Tibshirani et al. [191], defining the *gap index*

$$\gamma_{TWH}(p) = E(\mathcal{L}^*(p, \mathbb{T}^{\sharp})) - \mathcal{L}^*(p, T)$$

where the $\mathcal{L}^*(p, T)$ denotes the optimal value of the optimized cost with $p$ clusters for a training set $T$. The notation $\mathbb{T}^{\sharp}$ represent a random training set, with same size and dimension as $T$, generated using an unclustered probability distribution used as a reference. In Tibshirani et al. [191], this distribution is taken as uniform (over the smallest hypercube containing the observed data), or uniform on the coefficients of a principal component decomposition of the data (see chapter 20). The expectation $E(\mathcal{L}^*(p, \mathbb{T}^{\sharp}))$ is computed by Monte-Carlo simulation, by sampling many realizations of the training set $\mathbb{T}$, running the clustering algorithm for each of them and averaging the optimal costs.

One can expect $\mathcal{L}^*(p, T)$ (for observed data) to decrease much faster (when adding a cluster) than its expectation for homogeneous data when $p < p_0$, and the decrease of

both terms to be comparable when $p \geq p_0$. So the number of clusters can in principle be estimated by detecting an elbow in the graph of $\gamma_{TWH}(p)$ as a function of $p$. The procedure suggested in Tibshirani et al. [191] in order to detect this elbow if to look for the first index $p$ such that

$$\gamma_{TWH}(p+1) \leq \gamma_{TWH}(p) + \sigma(p+1)$$

where $\sigma(p+1)$ is the standard deviation of $\mathcal{L}^*(p+1, \mathbb{T}^\sharp)$ for homogeneous data, also estimated via Monte-Carlo simulation.

Figures figs. 19.5 to 19.7 provide a comparative illustration of some of these indexes.

## 19.7  Bayesian Clustering

### 19.7.1  Introduction

We have seen an example of model-based clustering with mixtures of Gaussian distributions. The main parameters in this model were the number of classes, $p$, and the probabilities $\alpha_j$ associated to each cluster, and the parameter of the conditional distribution (e.g., $\mathcal{N}(c_j, \sigma^2 \mathrm{Id}_{\mathbb{R}^d})$) of $X$ conditionally to being in the $j$th cluster. In the approach we described, these parameters were estimated from data using maximum likelihood (through the EM algorithm) and probabilities $f_Z(j|x)$ were then estimated in order to compute the most likely clustering. We interpreted $f_Z(j|x)$ as the conditional probability $P(Z = z|X = x)$, where $Z \in \{1, \ldots, p\}$ represents the group variable. The natural generative order is $Z \rightarrow X$: first decide to which group the observation belongs to, then sample the value of $X$ conditional to this group. Clustering is in this case reversing the order, i.e., computing the posterior distribution of $Z$ given $X$.

In a Bayesian approach, the parameters $p, \underline{\alpha}, \underline{c}$ and $\sigma^2$ are also considered as random variables, so that (letting $\underline{\theta}$ denote the vector formed by these parameters), the generative random sequence becomes $\underline{\theta} \rightarrow Z \rightarrow X$. Importantly, $\underline{\theta}$ is assumed to be generated once for all, even if several samples of $X$ are observed, yielding the generative sequence for an $N$-sample,

$$\underline{\theta} \rightarrow (Z_1, \ldots, Z_N) \rightarrow (X_1, \ldots, X_N).$$

We use below underlined letters to denote configurations of points, $\underline{Z} = (Z_1, \ldots, Z_N)$, $\underline{X} = (X_1, \ldots, X_N)$, etc. We also use capital letters or boldface letters (for Greek symbols) to differentiate random variable from realizations.

Clusters are still evaluated based on the conditional distribution of $\underline{Z}$ given $\underline{X}$, but this distribution must be evaluated by averaging the conditional distribution of

Figure 19.5: Comparison of cluster indices for Gaussian clusters. First row: original data and ground truth. Second panel: plots of four indices as functions of $p$ (Elbow; Caliński and Harabasz; silhouette; Sugar and James)

Figure 19.6:  Comparison of cluster indices for Gaussian clusters.  First row: original data and ground truth.  Second panel: plots of four indices as functions of $p$ (Elbow; Caliński and Harabasz; silhouette; Sugar and James).

Figure 19.7: Comparison of cluster indices for Gaussian clusters. First row: original data and ground truth. Second panel: plots of four indices as functions of $p$ (Elbow; Caliński and Harabasz; silhouette; Sugar and James).

$\underline{Z}$ and $\underline{\theta}$ given $\underline{X}$ with respect to $\underline{\theta}$, formally[4],

$$P(\underline{z}|\underline{x}) = \int P(\underline{z},\underline{\theta}|\underline{x})P(\underline{\theta})d\underline{\theta}$$

$$\propto \int \prod_{k=1}^{N} P(x_k|z_k,\underline{\theta})P(z_k|\underline{\theta})P(\underline{\theta})d\underline{\theta}.$$

In this expression, $P(\underline{\theta})d\underline{\theta}$ implies an integration with respect to the prior distribution of the parameters. This distribution is part of the design of the method, but one usually chooses it so that it leads to simple computations, using so-called *conjugate priors*, which are such that posterior distributions belong to the same parametric family as the prior. For example, the conjugate prior for the mean of a Gaussian distribution (such as $c_i$ in our model) is also a Gaussian distribution. The conjugate prior for a scalar variance is the inverse gamma distribution, with density

$$\frac{v^u}{\Gamma(u)}s^{-u-1}\exp(-v/s)$$

for some parameters $u,v$. A conjugate prior for the class probabilities $\underline{\alpha} = (\alpha_1,\dots,\alpha_p)$ is the Dirichlet distribution, with density

$$D(\alpha_1,\dots,\alpha_p) = \frac{\Gamma(a_1+\cdots+a_p)}{\Gamma(a_1)\cdots\Gamma(a_p)}\prod_{j=1}^{p}\alpha_j^{a_j-1}$$

on the simplex

$$\mathfrak{S}_p = \{(\alpha_1,\dots,\alpha_p)\in\mathbb{R}^p : \alpha_i \geq 0, \alpha_1+\cdots+\alpha_p = 1\}.$$

Note that these conjugate priors have the same form (up to normalization) as the parametric model densities when considered as functions of the parameters.

### 19.7.2  Model with a bounded number of clusters

We first discuss the Bayesian approach assuming that the number of clusters is smaller than a fixed number, $p$. In this example, we assume that $c_1,\dots,c_p$ are modeled as independent Gaussian variables $\mathcal{N}(0,\tau^2\mathrm{Id}_{\mathbb{R}^d})$, $\sigma^2$ with an inverse gamma distribution with parameters $u$ and $v$ and $(\alpha_1,\dots,\alpha_p)$ using a Dirichlet distribution with parameters $(a,\dots,a)$.

---

[4]The symbol $\propto$ means "equal up to a multiplicative constant".

**Analytical example.** The joint probability density of $(\underline{X}, \underline{Z})$ and $\underline{\theta}$ is proportional to

$$(\sigma^2)^{-u-1} e^{-v/\sigma^2} e^{-\sum_{j=1}^{p} |c_j|^2/2\tau^2} \prod_{j=1}^{p} \alpha_j^{a-1} \prod_{k=1}^{N} \frac{e^{-|x_k - c_{z_k}|^2/2\sigma^2}}{(\sigma^2)^{d/2}} \prod_{k=1}^{N} \alpha_{z_k}$$

$$= (\sigma^2)^{-u-dN/2-1} \exp\left(-(v + \frac{1}{2}\sum_{k=1}^{N} |x_k - c_{z_k}|^2)/\sigma^2\right) \prod_{j=1}^{p} \alpha_j^{a+N_j-1}.$$

One can explicitly integrate this last expression with respect to $\sigma^2$ and $\underline{\alpha}$, using the expressions of the normalizing constants in the inverse gamma and Dirichlet distributions, yielding (after integration and ignoring constant terms)

$$\frac{\Gamma(a + N_1)\cdots\Gamma(a + N_p)}{(v + \frac{1}{2}\sum_{k=1}^{N} |x_k - c_{z_k}|^2)^{u+dN/2}} \exp\left(-\sum_{j=1}^{p} |c_j|^2/2\tau^2\right)$$

$$= \frac{\Gamma(a + N_1)\cdots\Gamma(a + N_p)}{(v + \frac{1}{2}S_w + \frac{1}{2}\sum_{j=1}^{p} N_j |c_j - \bar{x}_j|^2)^{u+dN/2}} \exp\left(-\sum_{j=1}^{p} |c_j|^2/2\tau^2\right)$$

where $S_w = \sum_{k=1}^{N} |x_k - \bar{x}_{z_k}|^2$ is the within group sum of squares. Note that this sum of squares depends on $\underline{x}$ and $\underline{z}$, and that $(N_1, \ldots, N_p)$, the group sizes, depend on $\underline{z}$.

Let us assume a "non-informative prior" on the centers, which corresponds to letting $\tau$ tend to infinity and neglecting the last exponential. The remaining expression can now be integrated with respect to $c_1, \ldots, c_p$ by making a change of variables $\mu_j = \sqrt{N_j/(2v + S_k)}(c_j - \bar{x}_j)$ and using the fact that

$$\int_{(\mathbb{R}^d)^p} \frac{dc_1 \ldots dc_p}{(v + \frac{1}{2}S_w + \frac{1}{2}\sum_{j=1}^{p} N_j |c_j - \bar{x}_j|^2)^{u+dN/2}} =$$

$$(2v + S_w)^{(p-N)d/2-u)} \prod_{j=1}^{p} N_j^{-d/2} \int_{(\mathbb{R}^d)^p} \frac{d\mu_1 \ldots d\mu_p}{(\frac{1}{2} + \frac{1}{2}\sum_{j=1}^{p} |\mu_j|^2)^{u+dN/2}}$$

and the final integral does not depend on $\underline{x}$ or $\underline{z}$. It follows from this that the conditional distribution of $\underline{Z}$ given $\underline{x}$ takes the form

$$P(\underline{z}|\underline{x}) = C(\underline{x}) \frac{\prod_{j=1}^{p} \Gamma(a + N_j)}{(2v + S_w)^{(N-p)d/2+u)} \prod_{j=1}^{p} N_j^{d/2}}$$

where $C(\underline{x})$ is a normalization constant ensuring that the right-hand side is a probability distribution over configurations $\underline{z} = (z_1, \ldots, z_N) \in \{1, \ldots, p\}^N$. In order to obtain

the most likely configuration for this posterior distribution, one should therefore minimize in $\underline{z}$ the function

$$((N-p)\frac{d}{2}+u)\log(2v+S_w)+\frac{d}{2}\sum_{j=1}^{p}\log N_j-\sum_{j=1}^{p}\log\Gamma(a+N_j).$$

This final optimization problem cannot be solved in closed form, but this can be performed numerically. One can simplify it a little by only keeping the main order terms in the last two sums (using Stirling formula for the Gamma function) and minimize

$$((N-p)\frac{d}{2}+u)\log(2v+S_w)-\sum_{j=1}^{p}(a+N_j)\log(a+N_j).$$

This expression has a nice interpretation, since the first term minimizes the within-group sum of squares, the same objective function as in K-means, and the second one is an entropy term that favors clusters with similar sizes.

**Monte-Carlo simulation.**   An alternative to this analytical approach is to use Monte-Carlo simulations to estimate some properties of the posterior distribution numerically. While they are often computationally demanding, Monte-Carlo methods are more flexible and can be used in situations when analytic computations are intractable. In order to sample from the distribution of $\underline{Z}$ given $\underline{x}$, it is actually easier to sample from the joint distribution of $(\underline{Z}, \underline{\theta})$ given $\underline{x}$, because this distribution has a simpler form. Of course, if the pair $(\underline{Z}, \underline{\theta})$ is sampled from the conditional distribution given $\underline{x}$, the first component, $\underline{Z}$ will follow the posterior distribution we are interested in.

In the context of the discussed example, this reduces to sampling from a distribution proportional to

$$(\sigma^2)^{-u-1}e^{-v/\sigma^2}e^{-\sum_{j=1}^{p}|c_j|^2/2\tau^2}\prod_{j=1}^{p}\alpha_j^{a-1}\prod_{k=1}^{N}\frac{e^{-|x_k-c_{z_k}|^2/2\sigma^2}}{(\sigma^2)^{d/2}}\prod_{k=1}^{N}\alpha_{z_k}. \tag{19.11}$$

Sampling from all these variables at once is not tractable, but it is easy to sample from them in sub-groups, conditionally to the rest of the variables. We can, for example, deduce from the expression above the following conditional distributions.

(i) Given $(\underline{\alpha}, \underline{c}, \underline{z})$, $\sigma^2$ follows an inverse gamma distribution with parameters $u + dN/2$ and $v + \frac{1}{2}\sum_{k=1}^{N}|x_k - c_{z_k}|^2$.

(ii) Given $(\underline{z}, \underline{z}, \sigma^2)$, $\underline{\alpha}$ follows a Dirichlet distribution with parameters $a+N_1, \ldots, a+N_p$.

(iii) Given $(\underline{z}, \sigma^2, \underline{\alpha})$, $c_1, \ldots, c_p$ are independent and follow a Gaussian distribution, respectively with mean $(1 + \sigma^2/(N_j \tau^2))^{-1} \bar{x}_j$ and variance $(N_j/\sigma^2 + 1/\tau^2)^{-1}$.

(iv) Given $(\sigma^2, \underline{\alpha}, \underline{c})$, $z_1, \ldots, z_N$ are independent and

$$P(z_k = j | \sigma^2, \underline{\alpha}, \underline{c}, \underline{x}) \propto \alpha_j e^{-|x_k - c_j|^2/2\sigma^2}.$$

---

**Algorithm 19.10 (Gibbs sampling for mixture of Gaussian (Bayesian case))**
(1) Initialize with variables $\underline{\alpha}, \underline{c}, \sigma$ and $\underline{z}$, for example generated according to the prior distribution.

(2) Loop a large number of times over the following steps.

(i) Simulate a new value of $\sigma^2$ according to an inverse gamma distribution with parameters $u + dN/2$ and $v + \frac{1}{2} \sum_{k=1}^{N} |x_k - c_{z_k}|^2$.

(ii) Simulate new values for $\alpha_1, \ldots, \alpha_p$ according to a Dirichlet distribution with parameters $a + N_1, \ldots, a + N_p$.

(iii) Simulate new values for $c_1, \ldots, c_p$ independently, sampling $c_i$ according to a Gaussian distribution with mean $(1 + \sigma^2/(N_j \tau^2))^{-1} \bar{x}_j$ and variance $(N_j/\sigma^2 + 1/\tau^2)^{-1}$.

(iv) Simulate new values of $z_1, \ldots, z_N$ independently such that

$$P(z_k = j | \sigma^2, \underline{\alpha}, \underline{c}, \underline{x}) \propto \alpha_j e^{-|x_k - c_j|^2/2\sigma^2}.$$

---

Note that this algorithm is only asymptotically providing a sample of the posterior distribution (it has to be stopped at some point, of course). Note also that, at each step, the labels $z_1, \ldots, z_N$ provide a random partition of the set $\{1, \ldots, N\}$, and this partition changes at every step.

To estimate one single partition out of this simulation, several strategies are possible. Using the simulation, one can estimate the probability $w_{kl}$ that $x_k$ and $x_l$ belong to the same cluster. This can be dome by averaging the number of times that $z_k = z_l$ was observed along the Gibbs sampling iterations (from which one usually excludes a few early "burn-in" iterations). These weights, $w_{kl}$ can then be used as similarity measures in a clustering algorithm.

Alternatively, one can average for each $k$, the values of the class center $c_{z_k}$ associated to $k$, still along the Gibbs sampling iterations. These average values can then be used as input of, say, a K-means algorithm to estimate final clusters.

**Mean-field approximation.**   We conclude this section with a variational Bayes approximation of the posterior distribution. We will make a mean-field approximation, in which all parameters and latent variables are independent, therefore approximating the distribution in (19.11) by a product distribution taking the form

$$g(\sigma^2,\underline{\alpha},\underline{c},\underline{z}) = g^{(\sigma^2)}(\sigma^2)g^{(\alpha)}(\underline{\alpha})\prod_{j=1}^{p} g_j^{(c)}(c_j)\prod_{k=1}^{N} g_k^{(z)}(z_k).$$

Here $\underline{c} = (c_1,\ldots,c_p)$, $\underline{z} = (z_1,\ldots,z_N)$ and $\underline{\alpha} = (\alpha_1,\ldots,\alpha_p)$. We have $\sigma^2 \in (0,+\infty)$, $\underline{c} \in (\mathbb{R}^d)^p$, $\underline{\alpha} \in \mathcal{S}$, the set of all non-negative $\alpha_1,\ldots,\alpha_p$ that sum to one, and $\underline{z} \in \{1,\ldots,p\}^N$ (so that $g_k^{(x)}$ is a p.m.f. on $\{1,\ldots,p\}$. We will use the discussion in section 16.3.3 and lemma 16.1, and use the notation introduced in that section to denote as $\langle\varphi\rangle$ the expectation a variable $\varphi$ of the variables above for the p.d.f. $g$.

The log-likelihood for a mixture of Gaussian takes the form (ignoring contant terms)

$$\ell(\sigma^2,\underline{\alpha},\underline{c},\underline{z}) = -(u+1)\log\sigma^2 - v\sigma^{-2} - \frac{1}{2\tau^2}\sum_{k=1}^{p}|c_j|^2 + \sum_{j=1}^{p}(a-1)\log\alpha_j$$

$$- \frac{Nd}{2}\log\sigma^2 - \frac{1}{2}\sigma^{-2}\sum_{k=1}^{N}|x_k - c_{z_k}|^2 + \sum_{k=1}^{N}\log\alpha_{z_k}$$

$$= -(u+1)\log\sigma^2 - v\sigma^{-2} - \frac{1}{2\tau^2}\sum_{k=1}^{p}|c_j|^2 + \sum_{k=1}^{p}(a-1)\log\alpha_j$$

$$- \frac{Nd}{2}\log\sigma^2 - \frac{1}{2}\sigma^{-2}\sum_{k=1}^{N}\sum_{j=1}^{p}|x_k - c_j|^2\mathbf{1}_{z_k=j} + \sum_{k=1}^{N}\sum_{j=1}^{p}\log\alpha_j\mathbf{1}_{z_j=k}$$

and can therefore be decomposed as a sum of products of functions of single variables, as assumed in section 16.3.3. Using lemma 16.1, we can identify each of the distributions composing $g$, namely:

- $g^{(\sigma^2)}$ is the p.d.f. of an inverse gamma with parameters $\tilde{u} = u + Nd/2$ and

$$\tilde{v} = v + \frac{1}{2}\sum_{k=1}^{N}\sum_{j=1}^{p}\left\langle|x_k - C_j|^2\right\rangle\langle Z_k = j\rangle.$$

- $g_j^{(c)}$ is the p.d.f. of a Gaussian, with parameters $\mathcal{N}(\tilde{m}_j,\tilde{\sigma}_j^2\mathrm{Id}_{\mathbb{R}^d})$, with, letting

$$\tilde{\zeta}(j) = \sum_{k=1}^{N}\langle Z_k = j\rangle = \sum_{k=1}^{N}g_k^{(z)}(j),$$

$\tilde{\sigma}_j^2 = \left(\frac{1}{\tau^2} + \langle\sigma^{-2}\rangle\tilde{\zeta}(j)\right)^{-1}$ and $\tilde{m}_i = \langle\sigma^{-2}\rangle\tilde{\sigma}_j^2\sum_{k=1}^N\langle Z_k = j\rangle x_k$.

- $g^{(\alpha)}$ of a Dirichlet distribution, with parameters $\tilde{a}_1,\ldots,\tilde{a}_k$, with $\tilde{a}_i = a + \tilde{\zeta}(j)$.

- Finally $g_k^{(z)}$ is a p.m.f. on $\{1,\ldots,p\}$ with

$$g_k^{(z)}(j) \propto \exp\left(-\frac{1}{2}\langle\sigma^{-2}\rangle\langle|x_k - C_j|^2\rangle + \langle\log\alpha_j\rangle\right).$$

To complete the consistency equations, it now suffices to evaluate the expectations in the formula above as functions of the other parameters. We leave to the reader the verification of the following statements.

- If $\sigma^2$ follows an inverse gamma distribution with parameters $\tilde{u}$ and $\tilde{v}$, then $\langle\sigma^{-2}\rangle = \tilde{u}/\tilde{v}$.

- If $C_j \sim \mathcal{N}(\tilde{m}_j, \tilde{\sigma}_j^2\mathrm{Id}_{\mathbb{R}^d})$, then $\langle|x_k - C_j|^2\rangle = |x_k - \tilde{m}_j|^2 + d\tilde{\sigma}_j^2$.

- If $\underline{\alpha}$ follows a Dirichlet distribution with parameters $\tilde{a}_1,\ldots,\tilde{a}_p$, then $\langle\log\alpha_j\rangle = \psi(\tilde{a}_j) - \psi(\tilde{a}_1 + \cdots + \tilde{a}_p)$ where $\psi$ is the *digamma* function (derivative of the logarithm of the gamma function).

Combining these facts with the expression of the mean-field parameters, we can now formulate a mean-field estimation algorithm for mixtures of Gaussian that iteratively applies the consistency equations.

---

**Algorithm 19.11 (Mean-field algorithm for mixtures of Gaussian)**
(1) : Input: training set $(x_1,\ldots,x_N)$, number of clusters $p$, prior parameters $u,v,\tau^2$ and $a$.

(2) Initialize variables $\tilde{\sigma}_1^2,\ldots,\tilde{\sigma}_p^2,\tilde{m}_1,\ldots,\tilde{m}_p,\tilde{a}_1,\ldots,\tilde{a}_p,\tilde{g}_k(j)$, $k = 1,\ldots,N$, $j = 1,\ldots,p$.

(3) Let $\tilde{\zeta}(j) = \sum_{k=1}^N\tilde{g}_k(j)$, $j = 1,\ldots,p$.

(4) Let

$$\tilde{\rho}^2 = \frac{1}{u + Nd/2}\left(v + \frac{1}{2}\sum_{k=1}^N\sum_{j=1}^p\tilde{g}_k(j)|x_k - \tilde{m}_j|^2 + \frac{d}{2}\sum_{j=1}^p\tilde{\sigma}_j^2\tilde{\zeta}(j)\right).$$

(5) For $j = 1,\ldots,p$, let $\tilde{\sigma}_i^2 = \left(\frac{1}{\tau^2} + \frac{\tilde{\zeta}(j)}{\tilde{\rho}^2}\right)^{-1}$ and $\tilde{m}_i = \frac{\tilde{\sigma}_j^2}{\tilde{\rho}^2}\sum_{k=1}^N\tilde{g}_k(j)x_k$.

(6) Let $\tilde{a}_i = a + \tilde{\zeta}(j)$, $j = 1,\ldots,p$.

(7) For $k = 1,\ldots,N$, $j = 1,\ldots,p$, let

$$\tilde{g}_k(j) \propto \exp\left(-\frac{1}{2\tilde{\rho}^2}\left(|x_k - \tilde{m}_j|^2 + d\tilde{\sigma}_j^2\right) + \psi(\tilde{a}_j)\right).$$

(8) Compare the updated variables with their previous values and stop if the difference is below a tolerance level. Otherwise, return to (3).

---

After convergence $g_k^{(z)}$ provides the mean-field approximation of the posterior probability of classes for observation $k$ and can be used to determine clusters.

### 19.7.3   Non-parametric priors

**The Polya urn**   In the previous model with $p$ clusters or less, the joint distribution of $Z_1, \ldots, Z_N$ is given by

$$\pi(z_1, \ldots, z_N) = \frac{\Gamma(pa)}{\Gamma(a)^p} \int_{\mathfrak{S}_p} \prod_{j=1}^{p} \alpha_j^{a+N_j-1} \, d\alpha = \frac{\Gamma(pa)}{\Gamma(pa+N)} \prod_{j=1}^{p} \frac{\Gamma(a+N_j)}{\Gamma(a)}.$$

Conditional to $z_1, \ldots, z_N$, the data model was completed by sampling $p$ sets of parameters, say, $\theta_1, \ldots, \theta_p$, each belonging to a parameter space $\Theta$ and following a prior probability distribution with density, say, $\psi$ and variables $X_1, \ldots, X_N$, where $X_k \in \mathcal{R}$ was drawn according to a law dependent on its cluster, that we will denote $\varphi(\cdot \,|\, \theta_{z_k})$. The complete likelihood of the data is now

$$L(\underline{z}, \underline{\theta}, \underline{x}) = \frac{\Gamma(pa)}{\Gamma(pa+N)} \prod_{j=1}^{p} \frac{\Gamma(a+N_j)}{\Gamma(a)} \prod_{j=1}^{p} \psi(\theta_j) \prod_{k=1}^{N} \varphi(x_k | \theta_{z_k}).$$

Note that the right-hand side does not change if one relabels the values of $z_1, \ldots, z_N$, i.e., if one replaces each $z_k$ by $s(z_k)$ where $s$ is a permutation of $\{1, \ldots, p\}$, creating a new configuration denoted $s \cdot \underline{z}$. Let $[\underline{z}]$ denote the equivalence class of $\underline{z}$, containing all $\underline{z}' = s \cdot \underline{z}, s \in \mathfrak{S}_N$: all the labelings in $[\underline{z}]$ provide the same partition of $\{1, \ldots, N\}$ and can therefore be identified. One defines a probability distribution $\bar{\pi}$ over these equivalence classes by letting

$$\bar{\pi}([\underline{z}]) = |[\underline{z}]| \frac{\Gamma(pa)}{\Gamma(pa+N)} \prod_{j=1}^{p} \frac{\Gamma(a+N_j)}{\Gamma(a)}.$$

The first term on the right-hand side is the number of elements in the equivalence class of $[\underline{z}]$. To compute it, let $p_0 = p_0(\underline{z})$ denote the number of different values taken by $z_1, \ldots, z_N$, i.e., the "true" number of clusters (ignoring the empty ones), which now is a function of $\underline{z}$. Let $A_1, \ldots, A_{p_0}$ denote the partition associated with $\underline{z}$. New labelings equivalent to $\underline{z}$ can be obtained by assigning any index $i_1 \in \{1, \ldots, p\}$

to elements of $A_1$, then any index $i_2 \neq i_1$ to elements of $A_2$, etc., so that there are $\|[\underline{z}]\| = p!/(p - p_0)!$ choices. We therefore find:

$$\bar{\pi}([\underline{z}]) = \frac{p!}{(p - p_0)!} \frac{\Gamma(pa)}{\Gamma(pa + N)} \prod_{j=1}^{p} \frac{\Gamma(a + N_j)}{\Gamma(a)}.$$

Letting $\lambda = pa$ and using the formula $\Gamma(x + 1) = x\Gamma(x)$, this can be rewritten as

$$\bar{\pi}([\underline{z}]) = \frac{p(p - 1) \cdots (p - p_0 + 1)}{\lambda(\lambda + 1) \ldots (\lambda + N - 1)} \prod_{j=1}^{p} \prod_{i=0}^{N_j - 1} (\lambda/p + i).$$

Now, the class $[\underline{z}]$ contains exactly one element $\hat{\underline{z}}$ with the following properties

- $\hat{z}_1 = 1$,
- $\hat{z}_k \leq \max(z_j, j < k) + 1$ for all $k > 1$.

This means that the $k$th label is either one of those already appearing in $(\hat{z}_1, \ldots, \hat{z}_{k-1})$ or the next integer in the enumeration. We will call such a $\hat{\underline{z}}$ admissible. If we assume that $\underline{z}$ is admissible in the expression of $\bar{\pi}$, we can write

$$\bar{\pi}([\underline{z}]) = \frac{\prod_{j=1}^{p_0} \left( \lambda(1 - j/p) \prod_{i=1}^{N_j - 1} (\lambda/p + i) \right)}{\lambda(\lambda + 1) \ldots (\lambda + N - 1)}.$$

If one takes the limit $p \to \infty$ in this expression, one still gets a probability distribution on admissible labelings, namely

$$\bar{\pi}([\underline{z}]) = \frac{\lambda^{p_0} \prod_{j=1}^{p_0} (N_j - 1)!}{\lambda(\lambda + 1) \ldots (\lambda + N - 1)}. \tag{19.12}$$

Recall that, in this equation, $p_0$ is a function of $\underline{z}$, equal, for admissible labelings, to the largest $j$ such that $N_j > 0$.

The probability $\bar{\pi}$ is generated by the following sampling scheme, called the Polya urn process simulating admissible labelings.

---

**Algorithm 19.12 (Polya Urn)**

1 Initialize $k = 1$, $z_1 = 1$, $j = 1$. Let $N_1 = 1$

2 At step $k$, assume that $z_1, \ldots, z_k$ have been generated, with associated number of clusters equal to $j$ and $N_1, \ldots, N_j$ elements per cluster. Generate $z_{k+1}$ such that

$$z_{k+1} = \begin{cases} i & \text{with probability} \quad \dfrac{N_i}{\lambda + k}, \text{ for } i = 1, \ldots, j \\ j + 1 & \text{with probability} \quad \dfrac{\lambda}{\lambda + k} \end{cases} \tag{19.13}$$

3 If $z_{k+1} = i \le j$, then replace $N_i$ by $N_i + 1$, $k$ by $k + 1$.

4 If $z_{k+1} = j + 1$, let $N_{j+1} = 1$, replace $j$ by $j + 1$ and $k$ by $k + 1$.

5 If $k < N$, return to step 2, otherwise, stop.

---

Using this prior, the complete model for the distribution of the observed data is

$$L(\underline{z}, \underline{\theta}, \underline{x}) = \frac{\lambda^{p_0} \prod_{j=1}^{p_0}(N_j - 1)!}{\lambda(\lambda + 1)\dots(\lambda + N - 1)} \prod_{j=1}^{p_0} \psi(\theta_j) \prod_{k=1}^{N} \varphi(x_k | \theta_{z_k})$$

Recall that, in this expression, $\underline{z}$ is restricted to the set of admissible labelings. We also note that admissible labelings are in one-to-one correspondence with the partitions of $\{1, \dots, N\}$, so that the latent variable $\underline{z}$ in this expression can also be interpreted as representing a random partition of this set.

**Dirichlet processes.**   As we will see later, the expression of the global likelihood and the Polya urn model will suffice for us to develop non-parametric clustering methods for a set of observations $x_1, \dots, x_N$. However, this model is also associated to an important class of random probability distributions (i.e., random variables taking values in some set of probability distributions) called Dirichlet processes for which we provide a brief description.

The distribution in (19.12) was obtained by passing to the limit from a model that first generates $p$ numbers $\alpha_1, \dots, \alpha_p$, then generates the labels $z_1, \dots, z_N \in \{1, \dots, p\}$ identified modulo relabeling. This distribution can also be defined directly, by first defining an infinity of positive numbers $(\alpha_j, j \ge 1)$ such that $\sum_{i=1}^{\infty} \alpha_i = 1$, followed by the generation of random labels $Z_1, \dots, Z_N$ such that $P(Z_k = j) = \alpha_j$, followed once again with an identification up to relabeling.

The distribution of $\underline{\alpha}$ that leads to the Polya urn is called the *stick breaking process*. This process is such that

$$\alpha_j = U_j \prod_{i=1}^{j-1}(1 - U_i)$$

where $U_1, U_2, \dots$ is a sequence of i.i.d. variables following a Beta$(1, \lambda)$ distribution, i.e., with p.d.f. $\lambda(1 - u)^{\lambda - 1}$ for $u \in [0, 1]$. The stick breaking interpretation comes from the way $\alpha_1, \alpha_2, \dots$ can be simulated: let $\alpha_1 \sim$ Beta$(1, \lambda)$; given $\alpha_1, \dots, \alpha_{j-1}$, let $\alpha_j = (1 - \alpha_1 - \dots - \alpha_{j-1})U_j$ where $U_j \sim$ Beta$(1, \lambda)$ and is independent from the past. Each step can be thought of as breaking the remaining length, $(1 - \alpha_1 - \dots - \alpha_{j-1})$, of an original stick of length 1 using a beta-distributed variable, $U_j$. This process leads to the distribution (19.12) over admissible distributions, i.e., if $\alpha$ is generated according to the stick breaking process, and $Z_1, \dots, Z_N$ are independent, each such

that $P(Z_k = j) = \alpha_j$, then the probability that $(Z_1, \ldots, Z_N)$ is identical, after relabeling, to the admissible configuration $z$ is given by (19.12). (We skip the proof of this result, which is not straightforward.)

Now, take a realization $\underline{\alpha} = (\alpha_1, \alpha_2, \ldots)$ of the stick-breaking process, and independent realizations $\underline{\eta} = (\eta_1, \eta_1, \ldots)$ drawn according to the p.d.f. $\psi$. Define

$$\rho = \sum_{j=1}^{\infty} \alpha_j \delta_{\eta_j}. \tag{19.14}$$

For any realization of $\alpha$ and of $\eta$, $\rho$ is a probability distribution on the parameter space $\Theta$ (in which one chooses $\eta_i$ with probability $\alpha_i$). Since $\alpha$ and $\eta$ are both random variables, this defines a random variable $\rho$ with values in the space of probability measures on $\Theta$.

This process has the following characteristic property. For any family $V_1, \ldots, V_k \subset \Theta$ forming a partition of that set, the random variable $(\rho(U_1), \ldots, \rho(U_k))$ follows a Dirichlet distribution with parameters

$$\left( \lambda \int_{U_1} \psi \, d\eta, \ldots, \lambda \int_{U_1} \psi \, d\eta \right).$$

This is the definition of a Dirichlet process with parameters $(\lambda, \psi)$, or, simply, with parameter $\lambda\psi$. Conversely, one can also show that any Dirichlet process can be decomposed as in (19.14) where $\alpha$ is a stick-breaking process and $\eta$ independent realizations of $\psi$.

**Monte-Carlo simulation.** The joint distribution of labels, parameters and observed variables can also be deduced from (19.12), with a joint p.d.f. given by

$$\frac{\lambda^{p_0 - 1} \prod_{j=1}^{p_0} (N_j - 1)!}{(\lambda + 1) \cdots (\lambda + N - 1)} \prod_{j=1}^{p_0} \psi(\eta_j) \prod_{k=1}^{N} \varphi(x_k | \eta_{z_k}). \tag{19.15}$$

The forward simulation of this distribution is a straightforward extension of Algorithm 19.12, namely:

---

**Algorithm 19.13**

1 Initialize $k = 1$, $z_1 = 1$, $j = 1$. Let $N_1 = 1$.

2 Sample $\eta_1 \sim \psi$ and $x_1 \sim \varphi(\cdot | \eta_1)$.

3 At step $k$, assume that $z_1, \ldots, z_k$ has been generated, with associated number of clusters equal to $j$ and $N_1, \ldots, N_j$ elements per cluster. Generate $z_{k+1}$ such that

$$z_{k+1} = \begin{cases} i & \text{with probability} \quad \dfrac{N_i}{\lambda + k}, \text{ for } i = 1, \ldots, j \\ j + 1 & \text{with probability} \quad \dfrac{\lambda}{\lambda + k} \end{cases}$$

4 If $z_{k+1} = i \leq j$, sample $x_{k+1} \sim \varphi(\cdot | \eta_i)$. Replace $N_i$ by $N_i + 1$, $k$ by $k + 1$.

5 If $z_{k+1} = j + 1$, let $N_{j+1} = 1$, sample $\eta_{j+1} \sim \psi$ and $x_{k+1} \sim \varphi(\cdot | \eta_{j+1})$. Replace $j$ by $j + 1$ and $k$ by $k + 1$.

6 If $k < N$, return to step 2, otherwise, stop.

---

This algorithm cannot be used, of course, to sample from the conditional distribution of $Z$ and $\eta$ given $X = x$, and Markov-chain Monte-Carlo must be used for this purpose. In order to describe how Gibbs sampling may be applied to this problem, we use the fact that, as previously remarked, using admissible labelings $z$ is equivalent to using partitions $\mathcal{A} = (A_1, \ldots, A_{p_0})$ of $\{1, \ldots, N\}$, and we will use the latter formalism to describe the algorithm. We will also use the notation $\eta_A$ to denote the parameter associated to $A \in \mathcal{A}$ so our new notation for the variables is $(\mathcal{A}, \underline{\eta})$ where $\mathcal{A}$ is a partition of $\{1, \ldots, N\}$ and $\underline{\eta}$ is a collection $(\eta_A, A \in \mathcal{A})$ with $\eta_A \in \Theta$. Given this, we want to sample from a conditional p.d.f.

$$\Phi(\mathcal{A}, \underline{\eta} | x) \propto \frac{\lambda^{|\mathcal{A}|-1} \prod_{A \in \mathcal{A}} (|A| - 1)!}{(\lambda + 1) \cdots (\lambda + N - 1)} \prod_{A \in \mathcal{A}} \psi(\eta_A) \prod_{k \in A} \varphi(x_k | \eta_A). \tag{19.16}$$

As an additional notation, given a partition $\mathcal{A}$ and an index $k \in \{1, \ldots, N\}$, we let $\mathcal{A}_k$ denote the set $A$ in $\mathcal{A}$ that contains $k$.

The following points are relevant for the design of the sampling algorithm.

(1) The conditional distribution of $\underline{\eta}$ given $\mathcal{A}$ and the training data is proportional to

$$\prod_{A \in \mathcal{A}} \left( \psi(\underline{\eta}_A) \prod_{k \in A} \varphi(x_k | \underline{\eta}_A) \right)$$

This shows that the parameters $\eta_A, A \in \mathcal{A}$ are independent of each other, with $\eta_A$ following a distribution proportional to

$$\eta \mapsto \psi(\eta) \prod_{k \in A_j} \varphi(x_k | \eta).$$

Sampling from this distribution generally offers no special difficulty, especially if the prior $\psi$ is conjugate to $\varphi$. Importantly, one does not need to sample exactly from $\eta_A$, and it is often more convenient to separate $\eta_A$ into several components (such as mean and variance for mixtures of Gaussian) and sample from them alternatively, creating another level of Gibbs sampling.

(2) We now consider the issue of updating $\mathcal{A}$. We will use for this purpose the formalism of Algorithm 12.2. In particular, for each $k \in \{1, \ldots, N\}$, we associate to

the variable $(\mathcal{A}, \underline{\eta})$ the pair $(\mathcal{A}^{(k)}, \underline{\eta}^{(k)})$, where $\mathcal{A}^{(k)}$ is the partition of $\{1, \ldots, N\} \setminus \{k\}$ formed by the sets $A^{(k)} = A \setminus \{k\}$ and $\eta_A^{(k)} = \eta_A$, unless $A = \{k\}$, in which case the set and the corresponding $\eta_A$ are dropped.

We can write $\Phi(\mathcal{A}, \underline{\eta}|\underline{x})$ in the form

$$\Phi(\mathcal{A}, \underline{\eta}|\underline{x}) \propto q(\mathcal{A}_k, \eta_{\mathcal{A}_k}) \varphi(x_k|\eta_{\mathcal{A}_k}) \frac{\lambda^{|\mathcal{A}^{(k)}|-1} \prod_{B \in \mathcal{A}^{(k)}} (|B|-1)!}{(\lambda+1) \cdots (\lambda+N-1)} \prod_{B \in \mathcal{A}^{(k)}} \psi(\eta_B) \prod_{l \in B} \varphi(x_l|\eta_B)$$

(19.17)

with

$$q(A, \theta) = \sum_{B \in \mathcal{A}^{(k)}} |B| \mathbf{1}_{A = B \cup \{k\}} + \lambda \psi(\theta) \mathbf{1}_{A = \{k\}}$$

Partitions $\mathcal{A}'$ that are consistent with $\mathcal{A}^{(k)}$ allocate $k$ to one of the clusters in $\mathcal{A}^{(k)}$ or create a new cluster with a new parameter $\eta_k'$. If one replaces $(\mathcal{A}, \underline{\eta})$ by $(\mathcal{A}', \underline{\eta}')$, only the first two terms in (19.17) will be affected, so that the conditional probability of $\mathcal{A}'$ given $\mathcal{A}^{(k)}$ is proportional to $q(\mathcal{A}_k', \eta_{\mathcal{A}_k'}) \varphi(x_k|\eta_{\mathcal{A}_k'})$ and given by

$$\begin{cases} \dfrac{|B| \varphi(x_k|\eta_B)}{\mathfrak{C}_1 + \lambda \mathfrak{C}_2} & \text{if } \mathcal{A}_k' = B \cup \{k\}, \eta_B' = \eta_B, B \in \mathcal{A}^{(k)} \\[4mm] \dfrac{\lambda \varphi(x_k|\eta_k') \psi(\eta_k')}{\mathfrak{C}_1 + \lambda \mathfrak{C}_2} & \text{if } \mathcal{A}_k' = \{k\}, \end{cases}$$

where

$$\mathfrak{C}_1 = \sum_{B \in \mathcal{A}^k} |B| \varphi(x_k|\eta_B) \quad \text{and} \quad \mathfrak{C}_2 = \int_{\Theta} \varphi(x_k|\theta) \psi(\theta) d\theta.$$

Concretely, this means that one first decides to allocate $k$ to a set $B$ in $\mathcal{A}^{(k)}$ with probability $|B| \varphi(x_k|\eta_B)/(\mathfrak{C}_1 + \lambda \mathfrak{C}_2)$ and to create a new set with probability $\lambda \mathfrak{C}_2/(\mathfrak{C}_1 + \lambda \mathfrak{C}_2)$. If a new set is created, then the associated parameter $\eta_{\{k\}}'$ is sampled according to the p.d.f. $\varphi(x_k|\theta) \psi(\theta/\mathfrak{C}_2$.

(3) However, sampling using this conditional probability requires the computation of the integral $\mathfrak{C}_2$, which can represent a significant computational burden, since this has to be done many times in a Gibbs sampling algorithm. A modification of this algorithm, introduced in Neal [141], avoids this computation by adding new auxiliary variables at each step of the computation. These variables are $m$ parameters $\eta_1^*, \ldots, \eta_m^* \in \Theta$ where $m$ is a fixed integer. To define the joint distribution of $\mathcal{A}, \underline{\eta}, \underline{\eta}^*$, one lets the marginal distribution of $(\mathcal{A}, \underline{\eta})$ be given by (19.16) and conditionally to $\mathcal{A}, \underline{\eta}$, let $\eta_1^*, \ldots, \eta_m^*$ be:

(*i*) independent with density $\psi$ if $|\mathcal{A}_k| > 1$;

(*ii*) such that $\eta_j^* = \eta_{\mathcal{A}_k}$ and the other $m-1$ starred parameters are independent with distribution $\psi$, where $j$ is randomly chosen in $\{1, \ldots, m\}$ if $\mathcal{A}_k = \{k\}$.

With this definition, the joint conditional distribution of $(\mathcal{A}, \underline{\eta}, \underline{\eta}^*)$ takes the form

$$\widehat{\Phi}(\mathcal{A}, \underline{\eta}, \underline{\eta}^* | \underline{x}) \propto \hat{q}(\mathcal{A}_k, \eta_{\mathcal{A}_k}, \underline{\eta}^*) \varphi(x_k | \eta_{\mathcal{A}_k})$$

$$\frac{\lambda^{|\mathcal{A}^{(k)}|-1} \prod_{B \in \mathcal{A}^{(k)}}(|B|-1)!}{(\lambda+1)\cdots(\lambda+N-1)} \prod_{B \in \mathcal{A}^{(k)}} \psi(\eta_B) \prod_{l \in B} \varphi(x_l | \eta_B) \quad (19.18)$$

with

$$\hat{q}(A, \theta, \eta_1^*, \ldots, \eta_m^*) = \sum_{B \in \mathcal{A}^{(k)}} |B| \mathbf{1}_{\theta = \eta_B, A = B \cup \{k\}} \prod_{j=1}^{m} \psi(\eta_j^*) + \frac{\lambda}{m} \sum_{j=1}^{m} \mathbf{1}_{\theta = \eta_j^*, A = \{k\}} \psi(\theta) \prod_{i=1, i \neq j}^{m} \psi(\eta_i^*).$$

Note that $\widehat{\Phi}$ depends on $k$, so that the definition of the auxiliary variables will change at each step of Gibbs sampling. The conditional distribution, for $\widehat{\Phi}$, of $\mathcal{A}', \underline{\eta}'$ given $\mathcal{A}^{(k)}, \underline{\eta}^{(k)}, \underline{\eta}^*$ is such that

- $\mathcal{A}'_k = B \cup \{k\}$ and $\eta'_{\mathcal{A}'_k} = \eta_B$ with probability $|B| \varphi(x_k | \eta_B)/\mathfrak{C}$, for $B \in \mathcal{A}^{(k)}$.
- $\mathcal{A}'_k = \{k\}$ and $\eta_{\mathcal{A}'_k} = \eta_j^*$ with probability $(\lambda/m) \varphi(x_k | \eta_j^*)/\mathfrak{C}$, $j = 1, \ldots, m$.

The constant $\mathfrak{C}$ is given by

$$\mathfrak{C} = \sum_{B \in \mathcal{A}^k} |B| \varphi(x_k | \eta_B) + \frac{\lambda}{m} \sum_{j=1}^{m} \varphi(x_k | \eta_j^*)$$

and is therefore easy to compute.

We can now summarize this discussion with Neal's version of the Gibbs sampling algorithm.

---

**Algorithm 19.14 (Neal)**
Initialize the algorithm with some arbitrary partition and parameters $(\mathcal{A}, \underline{\eta})$ (for example, generated using the Dirichlet prior). Use the same notation to denote these variables at the end of the previous iteration of the algorithm. The next iteration is then run as follows.

(1) For $k = 1, \ldots, N$, reallocate $k$ to a cluster as follows.

   (i) Form the new family of sets $\mathcal{A}^{(k)}$ and labels $\underline{\eta}^{(k)}$ by removing $k$ from the partition $\mathcal{A}$.

   (ii) If $|\mathcal{A}_k| > 1$, generate $m$ variables $\eta_1^*, \ldots, \eta_m^*$ according to $\psi$. If $\mathcal{A}_k = \{k\}$, generate only $m-1$ such variables and let the last one be equal to $\eta_{\mathcal{A}_k}$.

(iii) Allocate $k$ to a new cluster $A'$ with parameter $\eta'_{A'}$ according to probabilities proportional to

$$\begin{cases} |B|\varphi(x_k|\eta_B^{(k)}) \text{ if } A' = B \cup \{k\} \text{ and } \eta'_{A'} = \eta_B^{(k)} \\ \dfrac{\lambda}{m}\varphi(x_k|\eta_j^*) \text{ if } A = \{k\} \text{ and } \eta'_{A'} = \eta_j^*, j = 1,\ldots,m \end{cases}$$

(2) For $A \in \mathcal{A}$, update $\eta_A, A \in \mathcal{A}$ according to the distribution proportional to

$$\psi(\eta)\prod_{k \in A}\varphi(x_k|\eta)$$

either directly, or via one step of Gibbs sampling visiting each of the variables that constitute $\eta_A$.

(3) Loop a sufficient number of times over the previous two steps.

---

After running this algorithm, the set of clusters should be finalized by using statistics computed along the simulation, as discussed after Algorithm 19.10.

**Full example: Mixture of Gaussian.** To conclude this section, we summarize the Monte-Carlo sampling algorithm for mixtures of Gaussian using a non-parametric Bayesian prior. Here, $\eta \in \Theta$ is the center $c \in \mathbb{R}^d$, with prior distribution $\psi = \mathcal{N}(0, \tau^2 \mathrm{Id}_{\mathbb{R}^d})$. The previous algorithm must be modified because an additional parameter $\sigma^2$ is shared by all classes, with prior given by an inverse gamma distribution with parameters $u$ and $v$. The conditional distribution of the data is $\varphi(x|c, \sigma) \sim \mathcal{N}(c, \sigma^2 \mathrm{Id}_{\mathbb{R}^d})$.

---

**Algorithm 19.15 (Gibbs sampling for non-parametric mixture of Gaussian)**
(1) Initialize the algorithm with some arbitrary partition and parameters $(\mathcal{A}, \underline{\eta})$.

(2) For $k = 1,\ldots,N$, reallocate $k$ to a cluster as follows.

(i) Form the new family of sets $\mathcal{A}^{(k)}$ and labels $\eta^{(k)}$ by removing $k$ from the partition $\mathcal{A}$.

(ii) If $|\mathcal{A}_k| > 1$, generate $m$ variables $c_i^*, i = 1,\ldots,m$ independently with $c_i^* \sim \mathcal{N}(0, \tau^2 \mathrm{Id}_{\mathbb{R}^d})$. If $\mathcal{A}_k = \{k\}$, generate only $m - 1$ such pairs of variables and let the last one be equal to $c_{\mathcal{A}_k}$.

(iii) Allocate $k$ to a new cluster $A'$ with parameter $c'_{A'}$ according to probabilities proportional to

$$\begin{cases} |B|\exp\Big(-\dfrac{|x_k - c_B^{(k)}|}{2\sigma^2}\Big) \text{ if } A' = B \cup \{k\} \text{ and } c'_{A'} = c_B^{(k)} \\ \dfrac{\lambda}{m}\exp\Big(-\dfrac{|x_k - c_B^*|}{2\sigma^2}\Big) \text{ if } A = \{k\} \text{ and } c'_{A'} = c_j^*, j = 1,\ldots,m. \end{cases}$$

(3) Simulate a new value of $\sigma^2$ according to an inverse gamma distribution with parameters $u + dN/2$ and $v + \frac{1}{2} \sum_{k=1}^{N} |x_k - c_{\mathcal{A}_k}|^2$.

(4) Simulate new values for $c_A, A \in \mathcal{A}$ independently, sampling $c_A$ according to a Gaussian distribution with mean $(1 + \sigma^2/(N_j \tau^2))^{-1} \bar{x}_A$ and variance $(|A|/\sigma^2 + 1/\tau^2)^{-1}$, where

$$\bar{x}_A = \frac{1}{|A|} \sum_{k \in A} x_k.$$

# Chapter 20

# Dimension Reduction and Factor Analysis

## 20.1 Principal component analysis

### 20.1.1 General Framework

Factor analysis aims at representing potentially high-dimensional data as functions of a (generally) small number of "factors," with a representation taking the general form

$$X = \Phi(Y, \theta) + \text{residual}, \tag{20.1}$$

where $X$ is the observation, $Y$ provide the factors and $\Phi$ is a function parametrized by $\theta$. A factor analysis model must therefore specify $\Phi$ (often, a linear function of $Y$), add hypotheses on $Y$ (such as its dimension, or properties of its distribution) and on the residuals. The transformation $\Phi$ is estimated from training data, but, ideally, the method should also provide an algorithm that infers $Y$ from a new observation of $X$. Most of the time, $Y$ is small dimensional so that the model also implies a reduction of dimension.

We start our discussion with principal component analysis (or PCA). This methods can be characterized in multiple ways, and we introducing through the angle of data approximation. In the following, the random variable $X$ takes values in a finite- or infinite-dimensional inner-product space $H$. We will denote, as usual, by $\langle . , . \rangle_H$ the product in this space.

Assume that $N$ independent realization of $X$, denoted $x_1, \ldots, x_N$, are observed, forming our training set $T$. Our goal is to obtain a small-dimensional representation of these data, while loosing a minimal amount of relevant information. PCA, is the simplest and most commonly used approach developed for this purpose.

If $V$ is a finite-dimensional subspace of $H$, we denote by $P_V(y)$ the orthogonal projection of $y \in H$ on $V$, i.e., the element $\xi \in V$ such that $\|y - \xi\|_H^2$ is minimal

497

(see section 6.4). Recall that this orthogonal projection if characterized by the two properties: (i) $P_V(y) \in V$ and (ii) $(y - P_V(y)) \perp V$.

Given a target dimension $p$, PCA determines a $p$-dimensional subspace of $H$, say, $V$ and a point $c \in H$, such that, letting

$$R_k = x_k - c - P_V(x_k - c)$$

for $k = 1, \ldots, N$, the residual sum of squares

$$S = \sum_{k=1}^{N} \|R_k\|_H^2 \tag{20.2}$$

is as small as possible.

An optimal choice for $c$ is $c = \bar{x} = \sum_{k=1}^{N} x_k / N$. Indeed, using the linearity of the orthogonal projection, we have

$$S = \sum_{k=1}^{N} \|x_k - P_V(x_k) - (c - P_V(c))\|_H^2$$

$$= \sum_{k=1}^{N} \|x_k - P_V(x_k) - (\bar{x} - P_V(\bar{x}))\|_H^2 + N\|\bar{x} - P_V(\bar{x}) - (c - P_V(c))\|_H^2.$$

Given this, there would be no loss of generality in assuming that all $x_k$'s have been replaced by $x_k - \bar{x}$ and taking $c = 0$. While this is often done in the literature, there are some advantages (especially when discussing kernel methods) in keeping the average explicit in the notation, as we will continue to do.

Introducing an orthonormal basis $(e_1, \ldots, e_p)$ of $V$, one has

$$P_V(x_k - \bar{x}) = \sum_{i=1}^{p} \rho_k(i) e_i$$

with $\rho_{ki} = \langle x_k - \bar{x}, e_i \rangle_H$. One can then reformulate the problem in terms of $(e_1, \ldots, e_p)$, which must minimize

$$\begin{aligned}
S &= \sum_{k=1}^{N} \|x_k - \bar{x} - \sum_{i=1}^{p} \langle x_k - \bar{x}, e_i \rangle e_i\|_H^2 \\
&= \sum_{k=1}^{N} \|x_k - \bar{x}\|_H^2 - \sum_{i=1}^{p} \sum_{k=1}^{N} \langle x_k - \bar{x}, e_i \rangle_H^2.
\end{aligned}$$

For $u, v \in H$, define

$$\langle u, v \rangle_T = \frac{1}{N} \sum_{k=1}^{N} \langle x_k - \bar{x}, u \rangle_H \langle x_k - \bar{x}, v \rangle_H$$

and $\|u\|_T = \langle u, u \rangle_T^{1/2}$ (the index $T$ refers to the fact that this norm is associated with the training set). This provides a new quadratic form on $H$. The formula above shows that minimizing $S$ is equivalent to maximizing

$$\sum_{i=1}^{p} \|e_i\|_T^2$$

subject to the constraint that $(e_1, \ldots, e_p)$ is orthonormal in $H$.

Let us consider a slightly more general problem. If $H$ is a separable Hilbert space[1] and $\mu$ is a square-integrable probability measure on $H$, such that

$$\int_H \|x\|_H^2 \, d\mu(x) < \infty,$$

one can define $m = \int_H x \, d\mu(x)$ and $\sigma_\mu^2 = \int_H \|x - m\|_H^2 \, d\mu$. One can then define the *covariance bilinear form*

$$\Gamma_\mu(u, v) = \int_H \langle u, x - m \rangle_H \langle v, x - m \rangle_H \, d\mu(x),$$

which satisfies $\Gamma_\mu(u, v) \leq \sigma_\mu^2 \|u\|_H \|v\|_H$.

With this notation, we have

$$\langle u, v \rangle_T = \Gamma_{\hat{\mu}_T}(u, v),$$

where $\hat{\mu}_T = (1/N) \sum_{k=1}^{N} \delta_{x_k}$ is the empirical measure (and in that case $m = \bar{x}$). We can therefore generalize the PCA problem by considering the maximization of

$$\sum_{k=1}^{p} \Gamma_\mu(e_k, e_k) \tag{20.3}$$

over all orthonormal families $(e_1, \ldots, e_p)$ in $H$.

When $\mu$ is square integrable, the associated operator, $A_\mu$ defined by

$$\langle u, A_\mu v \rangle_H = \Gamma_\mu(u, v) \tag{20.4}$$

---

[1]A Hilbert space is an inner-product space which is complete for its norm. A separable Hilbert space must have a dense countable subset, which, in particular, implies that it has orthonormal bases.

for all $u, v \in H$, is a *Hilbert-Schmidt operator* [205]. Such an operator can, in particular, be diagonalized in an orthonormal basis of $H$, i.e., there exists an orthonormal basis, $(f_1, f_2, \dots)$ of $H$ such that $A_\mu f_i = \lambda_i^2 f_i$ for a non-increasing sequence of eigenvalues (with $\lambda_1 \geq \lambda_2 \geq \cdots \geq 0$) such that

$$\sigma_\mu^2 = \sum_{k=1}^\infty \lambda_i^2.$$

The main statement of the following result is in finite dimensions, a simple application of corollary 2.4. We here give a direct proof that also works in infinite dimensions.

**Theorem 20.1** *Let $(f_1, f_2, \dots)$ be an orthonormal basis of eigenvectors of $A_\mu$ with associated eigenvalues $\lambda_1^2 \geq \lambda_2^2 \geq \cdots \geq 0$. Then an orthonormal family $(e_1, \dots, e_p)$ in $H$ maximizes* (20.3) *if and only if,*

$$\mathrm{span}(f_j : \lambda_j^2 > \lambda_p^2) \subset \mathrm{span}(e_1, \dots, e_p) \subset \mathrm{span}(f_j : \lambda_j^2 \geq \lambda_p^2). \tag{20.5}$$

*In particular $f_1, \dots, f_p$ always provide a solution and $\mathrm{span}(e_1, \dots, e_p) = \mathrm{span}(f_1, \dots, f_p)$ for any other solution as soon as $\lambda_p^2 > \lambda_{p+1}^2$.*

**Definition 20.2** *When $\mu = \hat{\mu}_T$, the vectors $(f_1, \dots, f_p)$ are called (with some abuse when eigenvalues coincide) the first $p$ principal components of the training set $(x_1, \dots, x_N)$.*

PROOF  If $(e_1, \dots, e_p)$ is an orthonormal family in $H$, let

$$F(e_1, \dots, e_p) = \sum_{k=1}^p \Gamma_\mu(e_k, e_k).$$

Note that $F(f_1, \dots, f_p) = \lambda_1^2 + \cdots + \lambda_p^2$. Write $e_k = \sum_{j=1}^\infty \alpha_k^{(j)} f_j$ (so that $\alpha_k^{(j)} = \langle f_j, e_k \rangle_H$). These coefficients satisfy $\sum_{j=1}^\infty \alpha_k^{(j)} \alpha_l^{(j)} = 1$ if $k = l$ and 0 otherwise. Then

$$\Gamma_\mu(e_k, e_k) = \sum_{j=1}^\infty \lambda_j^2 (\alpha_k^{(j)})^2.$$

We have

$$F(e_1,\ldots,e_p) = \sum_{k=1}^{p} \sum_{j=1}^{\infty} \lambda_j^2 (\alpha_k^{(j)})^2$$

$$= \sum_{k=1}^{p} \sum_{j=1}^{p} \lambda_j^2 (\alpha_k^{(j)})^2 + \sum_{k=1}^{p} \sum_{j=p+1}^{\infty} \lambda_j^2 (\alpha_k^{(j)})^2$$

$$\leq \sum_{k=1}^{p} \sum_{j=1}^{p} \lambda_j^2 (\alpha_k^{(j)})^2 + \sum_{k=1}^{p} \sum_{j=p+1}^{\infty} \lambda_{p+1}^2 (\alpha_k^{(j)})^2$$

$$= \sum_{j=1}^{p} (\lambda_j^2 - \lambda_{p+1}^2) \sum_{k=1}^{p} (\alpha_k^{(j)})^2 + p\lambda_{p+1}^2.$$

Let $P$ denote the orthogonal projection operator from $H$ to $\mathrm{span}(e_1,\ldots,e_p)$. We have, for any $h \in H$, $\|Ph\|_H^2 \leq \|h\|_H^2$ with equality if and only if $h \in \mathrm{span}(e_1,\ldots,e_p)$. Applying this to $h = f_j$, with $P(f_j) = \sum_{k=1}^{p} \alpha_k^{(j)} e_k$, we get $\sum_{k=1}^{p}(\alpha_k^{(j)})^2 \leq 1$ with equality if and only if $f_j \in \mathrm{span}(e_1,\ldots,e_p)$.

As a consequence, the previous upper bound on $F(e_1,\ldots,e_p)$ implies

$$F(e_1,\ldots,e_p) \leq \sum_{j=1}^{p} \lambda_j^2.$$

This upper bound is attained at $(e_1,\ldots,e_p) = (f_1,\ldots,f_p)$, which is therefore a maximizer. Also, inspecting the argument above, we see that $F(e_1,\ldots,e_p) < \lambda_1^2 + \cdots + \lambda_p^2$ unless

(a) for all $k \leq p$ and $j \geq p+1$: $\alpha_k^{(j)} = 0$ if $\lambda_j^2 > \lambda_{p+1}^2$, and

(b) for all $j \leq p$: $\sum_{k=1}^{p}(\alpha_k^{(j)})^2 = 1$ unless $\lambda_j^2 = \lambda_{p+1}^2$.

Condition (a) implies that $\mathrm{span}(e_1,\ldots,e_p) \subset \mathrm{span}(f_j : \lambda_j^2 \leq \lambda_{p+1}^2)$. If $\lambda_p^2 = \lambda_{p+1}^2$, the inclusion $\mathrm{span}(e_1,\ldots,e_p) \subset \mathrm{span}(f_j : \lambda_j^2 \leq \lambda_p^2)$ therefore holds. If $\lambda_p^2 < \lambda_{p+1}^2$, condition (b) requires $\sum_{k=1}^{p}(\alpha_k^{(j)})^2 = 1$ for all $j \leq p$, which implies $f_j \in \mathrm{span}(e_1,\ldots,e_p)$ for $j \leq p$, so that $\mathrm{span}(e_1,\ldots,e_p) = \mathrm{span}(f_1,\ldots,f_p)$ and the inclusion also hold.

Condition (b) always requires $\sum_{k=1}^{p}(\alpha_k^{(j)})^2 = 1$, hence $f_j \in \mathrm{span}(f_1,\ldots,f_p)$, when $\lambda_j < \lambda_p$, showing that $\mathrm{span}(f_j : \lambda_j^2 < \lambda_p^2) \subset \mathrm{span}(e_1,\ldots,e_p)$. Equation (20.5) therefore always holds for $(e_1,\ldots,e_p)$ such that $F(e_1,\ldots,e_p) = \lambda_1^2 + \cdots + \lambda_p^2$. Furthermore, conditions (a) and (b) always hold for any orthonormal family that satisfy (20.5), showing that any such solution is optimal. ∎

Notice that the optimal $S$ in (20.2) is such that

$$S = N \sum_{i>p} \lambda_i^2.$$

**Remark 20.3** The interest of discussing PCA associated with a covariance operator for a square integrable measure (in which case it is often called a Karhunen-Loeve (KL) expansion) is that this setting is often important when discussing infinite-dimensional random processes (such as Gaussian random fields). Moreover, these operators quite naturally provide asymptotic versions of sample-based PCA. Interesting issues, that are part of *functional data analysis* [158], address the design of proper estimation procedures to obtain converging estimators of KL expansions based on finite samples for stochastic processes in infinite-dimensional spaces.  ◆

### 20.1.2   Computation of the principal components

**Small dimension.**   Assume that $H$ has finite dimension, $d$, i.e., $H = \mathbb{R}^d$, and represent $x_1, \ldots, x_N \in \mathbb{R}^d$ as column vectors. Let the inner product on $H$ be associated to a positive-definite symmetric matrix $Q$:

$$\langle u, v \rangle_H = u^T Q v.$$

Introduce the covariance matrix of the data

$$\Sigma_T = \frac{1}{N} \sum_{k=1}^{N} (x_k - \overline{x})(x_k - \overline{x})^T,$$

Write $A_T = A_{\hat{\mu}_T}$, for short, in (20.4). We have:

$$\langle u, A_T v \rangle_H = \frac{1}{N} \sum_{k=1}^{N} (u^T Q(x_k - \overline{x}))(v^T Q(x_k - \overline{x}))$$

$$= \frac{1}{N} \sum_{k=1}^{N} u^T Q(x_k - \overline{x})(x_k - \overline{x})^T Q v$$

$$= \langle u, \Sigma_T Q v \rangle_H,$$

so that $A_T = \Sigma_T Q$.

The eigenvectors, $f$, of $A_T$ are such that $Q^{1/2} f$ are eigenvectors of the symmetric matrix $Q^{1/2} \Sigma_T Q^{1/2}$, which shows that they form an orthogonal system in $H$, which will be orthonormal if the eigenvectors are normalized so that $f^T Q f = 1$. Equivalently, they solve the generalized eigenvalue problem $Q \Sigma_T Q f = \lambda^2 Q f$, which may be preferred numerically to diagonalizing the non-symmetric matrix $\Sigma_T Q$.

**Remark 20.4** Sometimes, the metric is specified by giving $Q^{-1}$ instead of $Q$ (or $Q^{-1}$ is easy to compute). Then, one can directly solve the generalized eigenvalue problem $\Sigma_T \tilde{f} = \lambda^2 Q^{-1} \tilde{f}$ and set $f = Q^{-1}\tilde{f}$. The normalization $f^T Q f = 1$ is then obtained by normalizing $\tilde{f}$ so that $\tilde{f}^T Q^{-1} \tilde{f} = 1$. ◆

**Remark 20.5** The "standard" version of PCA applies this computation using the Euclidean inner product, with $Q = \mathrm{Id}_{\mathbb{R}^d}$, and the principal components are the eigenvectors of the covariance matrix of $T$ associated with the largest eigenvalues. ◆

**Large dimension.** It often happens that the dimension of $H$ is much larger than the number of observations, $N$. In such a case, the previous approach is quite inefficient (especially when the dimension of $H$ is infinite!) and one should proceed as follows.

Returning to the original problem, one can remark that there is no loss of generality in assuming that $V$ is a subspace of $W := \mathrm{span}\{x_1 - \overline{x}, \dots, x_N - \overline{x}\}$. Indeed, letting $V' = P_W(V)$ (the projection of $V$ on $W$), we have, for $\xi \in W$,

$$
\begin{aligned}
\|\xi - P_V \xi\|_H^2 &= \|\xi\|_H^2 - 2\langle \xi, P_V \xi \rangle_H + \|P_V \xi\|_H^2 \\
&= \|\xi\|_H^2 - 2\langle P_W \xi, P_V \xi \rangle_H + \|P_V \xi\|_H^2 \\
&= \|\xi\|_H^2 - 2\langle \xi, P_W P_V \xi \rangle_H + \|P_V \xi\|_H^2 \\
&\geq \|\xi\|_H^2 - 2\langle \xi, P_W P_V x \rangle_H + \|P_W P_V \xi\|_H^2 \\
&= \|\xi - P_W P_V \xi\|_H^2 \\
&\geq \|\xi - P_{V'} \xi\|_H^2 .
\end{aligned}
$$

In this computation, we have used the facts that $P_W \xi = \xi$ (since $\xi \in W$), that $\|P_W P_V \xi\|_H \leq \|P_V \xi\|_H$, that $P_W P_V \xi \in V'$ and that $P_{V'}(\xi)$ is the best approximation of $\xi$ by an element of $V'$. This shows that (since $x_k - \overline{x} \in W$ for all $k$)

$$
\sum_{k=1}^N \|x_k - \overline{x} - P_V(x_k - \overline{x})\|_H^2 \geq \sum_{k=1}^N \|x_k - \overline{x} - P_{V'}(x_k - \overline{x})\|_H^2
$$

with $V'$ a subspace of $W$ of dimension less than $p$, proving the result. This computation also shows that no improvement in PCA can be obtained by looking for spaces of dimension $p \geq \dim(W)$ (with $\dim(W) \leq N - 1$ because the data is centered).

It therefore suffices to look for $f_1, \dots, f_p$ in the form

$$
f_i = \sum_{k=1}^N \alpha_k^{(i)} (x_k - \overline{x}).
$$

for some $\alpha_k^{(i)}$, $1 \leq k \leq N, 1 \leq i \leq p$.

With this notation, we have $\langle f_i, f_j \rangle_H = \sum_{k,l=1}^{N} \alpha_k^{(i)} \alpha_l^{(j)} \langle x_k - \overline{x}, x_l - \overline{x} \rangle_H$ and

$$\langle f_i, f_j \rangle_T = \frac{1}{N} \sum_{l=1}^{N} \langle f_i, x_l - \overline{x} \rangle_H \langle f_j, x_l - \overline{x} \rangle_H$$

$$= \frac{1}{N} \sum_{k,k'=1}^{N} \alpha_k^{(i)} \alpha_{k'}^{(j)} \sum_{l=1}^{N} \langle x_k - \overline{x}, x_l - \overline{x} \rangle_H \langle x_{k'} - \overline{x}, x_l - \overline{x} \rangle_H.$$

Let $S$ be the Gram matrix of the centered data, formed by the inner products $\langle x_k - \overline{x}, x_l - \overline{x} \rangle_H$, for $k, l = 1, \dots, N$. Let $\alpha^{(i)}$ be the column vector with coordinates $\alpha_k^{(i)}$, $k = 1, \dots, N$. We have $\langle f_i, f_j \rangle_H = (\alpha^{(i)})^T S \alpha^{(j)}$ and $\langle f_i, f_j \rangle_T = (\alpha^{(i)})^T S^2 \alpha^{(j)} / N$, which implies that, in this representation, the operator $A_T$ is given by $S/N$. Thus, the previous simultaneous orthogonalization problem can be solved in terms of the $\alpha$'s by diagonalizing $S$ and taking the first eigenvectors, normalized so that $(\alpha^{(i)})^T S \alpha^{(i)} = 1$. Let $\lambda_j^2$, $j = 1, \dots, N$ be the eigenvalues of $S/N$ (of which only the first $\min(d, N-1)$ may be non-zero). In this representation, the decomposition of the projection of $x_k$ on the PCA basis is given by

$$x_k = \sum_{j=1}^{p} \beta_k^{(j)} f_j$$

with

$$\beta_k^{(j)} = \langle x_k - \overline{x}, f_j \rangle_H = \sum_{l=1}^{N} \alpha_l^{(j)} \langle x_l - \overline{x}, x_k - \overline{x} \rangle_H = N \lambda_j^2 \alpha_k^{(j)}.$$

## 20.2  Kernel PCA

Since the previous computation only depended on the inner products $\langle x_k - \overline{x}, x_l - \overline{x} \rangle_H$, PCA can be performed in reproducing kernel Hilbert spaces, and the resulting method is called kernel PCA. In this framework, $X$ may take values in any set $\mathcal{R}$ with a representation $h : \mathcal{R} \to H$. The associated kernel, $K(x, x') = \langle h(x), h(x') \rangle_H$, provides a closed form expression of the inner products in terms of the original variables. The feature function itself is most of the time unnecessary.

The kernel version of PCA consists in replacing $x_k - \overline{x}$ with $h(x_k) - \overline{h}$ where $\overline{h}$ is the average feature. This leads to defining a "centered kernel:"

$$\begin{aligned} K_c(x, x') &= \langle h(x) - \overline{h}, h(x') - \overline{h} \rangle_H \\ &= \langle h(x), h(x') \rangle_H - \langle h(x) + h(x'), \overline{h} \rangle + \left\| \overline{h} \right\|_H^2 \\ &= K(x_k, x_l) - \frac{1}{N} \sum_{k=1}^{N} (K(x, x_k) + K(x', x_k)) + \frac{1}{N^2} \sum_{k,l=1}^{N} K(x_k, x_l). \end{aligned}$$

Then the Gram matrix in feature space is $S$ with $s_{kl} = K_c(x_k, x_l)$ and the computation described in the previous section can be applied. Note that, if one denotes, as usual $\mathcal{K} = \mathcal{K}(x_1, \ldots, x_N)$ the matrix formed by kernel evaluations $K(x_k, x_l)$, and if one lets $P = \mathrm{Id}_{\mathbb{R}^N} - \mathbb{1}_N \mathbb{1}_N / N$, then we have the simple matrix expression $S = P\mathcal{K}P$.

Letting $\alpha^{(1)}, \ldots, \alpha^{(p)} \in \mathbb{R}^N$ be the first $p$ eigenvectors of $S$, normalized so that $(\alpha^{(i)})^T S \alpha^{(i)} = 1$, the principal directions are vectors in feature space given by (using the notation in the previous section in which the $k$th coordinate of $\alpha^{(i)}$ is $\alpha_k^{(i)}$)

$$f_i = \sum_{k=1}^{N} \alpha_k^{(i)} (h(x_k) - \bar{h}),$$

and they are not computable when the features not known explicitly. However, a few geometric features associated with these directions can be characterized using the kernel only.

Consider the line in feature space $D_i = \{\bar{h} + \lambda f_i, \lambda \in \mathbb{R}\}$. Let $\Omega_i$ denote the points $x \in \mathcal{R}$ such that $h(x) \in D_i$. Then $x \in \Omega_i$ if and only if $h(x)$ coincides with its orthogonal projection on $D_i$, which is equivalent to

$$\langle h(x) - \bar{h}, f_i \rangle_H^2 = \left\| h(x) - \bar{h} \right\|_H^2,$$

which can be expressed with the kernel as

$$K_c(x, x) - \left( \sum_{k=1}^{N} \alpha_k^{(i)} K_c(x, x_k) \right)^2 = 0. \tag{20.6}$$

This provides a nonlinear equation in $x$. In particular, $\Omega_i$ is generally nonlinear, possibly with several connected components. Note that, by definition, the difference in (20.6) is always non-negative, so that a way to visualize $\Omega_i$ is to compute its sublevel sets, i.e., the set of all $x$ such that

$$K_c(x, x) - \left( \sum_{k=1}^{N} \alpha_k^{(i)} K_c(x, x_k) \right)^2 \leq \epsilon$$

for small $\epsilon$.

Similarly, the feature vector $h(x) - \bar{h}$ belongs to the space generated by the first $p$ components if and only if

$$\sum_{i=1}^{p} \langle h(x) - \bar{h}, f_i \rangle_H^2 = \left\| h(x) - \bar{h} \right\|_H^2$$

i.e.,

$$\sum_{i=1}^{p} \left( \sum_{k=1}^{N} \alpha_k^{(i)} K_c(x, x_k) \right)^2 = K_c(x, x).$$

One can also compute the finite-dimensional coordinates of $h(x)$ in the PCA basis, and this computation is easier. The representation is

$$x \mapsto (u_1(x), \ldots, u_p(x))$$

with

$$u_i = \langle h(x) - \bar{h}, f_i \rangle_H = \sum_{k=1}^{N} \alpha_k^{(i)} K_c(x, x_k).$$

This provides an explicit nonlinear transformation that maps each data point $x$ into a $p$-dimensional point. This representation allows one to easily exploit the reduction of dimension.

## 20.3   Statistical interpretation and probabilistic PCA

There is a simple probabilistic interpretation of linear PCA. Assume that $H = \mathbb{R}^d$ with the standard inner product and that $X$ is a centered random vector with covariance matrix $\Sigma$. Consider the problem that consists in finding a factor decomposition

$$X = \sum_{i=1}^{p} Y^{(i)} e_i + R$$

where $Y = (Y^{(1)}, \ldots, Y^{(p)})^T$ forms a $p$-dimensional centered vector, $e_1, \ldots, e_p$ is an orthonormal system, and $R$ is a random vector, independent of $Y$ and as small as possible, in the sense that $E(|R|^2)$ is minimal.

One can see that, in an optimal decomposition, one needs $R^T e_i = 0$ for all $i$, because one can always write

$$\sum_{i=1}^{p} Y^{(i)} e_i + R = \sum_{i=1}^{p} (Y^{(i)} + R^T e_i) e_i + R - \sum_{i=1}^{p} R^T e_i e_i.$$

If $R$ is centered, then so is $R - \sum_{i=1}^{p} R^T e_i e_i$ and the latter provides a better solution since $|R - \sum_{i=1}^{p} R^T e_i e_i| \leq |R|$. Also, there is no loss of generality in requiring that $(Y^{(1)}, \ldots, Y^{(p)})$ are uncorrelated, as this can always be obtained after a change of basis in $\text{span}(e_1, \ldots, e_p)$.

Assuming this, we can write

$$E(|X|^2) = \sum_{i=1}^{p} E((Y^{(i)})^2) + E(|R|^2)$$

with $Y^{(i)} = e_i^T X$. So, to minimize $E(|R|^2)$, one needs to maximize

$$\sum_{i=1}^{p} E((e_i^T X)^2)$$

which is equal to (letting $\Sigma$ be the covariance matrix of $X$)

$$\sum_{i=1}^{p} e_i^T \Sigma e_i.$$

The solution of this problem is given by the first $p$ eigenvectors of $\Sigma$. PCA (with a Euclidean metric) exactly applies this procedure, with $\Sigma$ replaced by the empirical covariance.

"Probabilistic PCA" is based on a slightly different statistical model in which it is assumed that $X$ can be decomposed as

$$X = \sum_{i=1}^{p} \lambda_i Y^{(i)} e_i + \sigma R,$$

where $R$ is a $d$ dimensional standard Gaussian vector and $Y = (Y^{(1)}, \ldots, Y^{(p)})^T$ a $p$-dimensional standard Gaussian vector, independent of $R$. The main difference with standard PCA is that the total variance of the residual, here $d\sigma^2$, is a model parameter and not a quantity to minimize.

In addition to $\sigma^2$, the model is parametrized by the coordinates of $e_1, \ldots, e_p$ and the values of $\lambda_1, \ldots, \lambda_p$. Introduce the $d \times p$ matrix

$$W = [\lambda_1 e_1, \ldots, \lambda_p e_p].$$

We can rewrite this model in the form

$$X = WY + \sigma^2 R$$

where the parameters are $W$ and $\sigma^2$, with the constraint that $W^T W$ is a diagonal matrix. As a linear combination of independent Gaussian random variables, $X$ is

Gaussian with covariance matrix $WW^T + \sigma^2\mathrm{Id}$. The log-likelihood of the observations $x_1, \ldots, x_N$ therefore is

$$L(W, \sigma) = -\frac{N}{2}\left(d\log 2\pi + \log\det(WW^T + \sigma^2\mathrm{Id}) + \mathrm{trace}((WW^T + \sigma^2\mathrm{Id})^{-1}\Sigma_T)\right) \quad (20.7)$$

where $\Sigma_Y$ is the empirical covariance matrix of $x_1, \ldots, x_N$. This function can be maximized explicitly in $W$ and $\sigma$, as stated in the following proposition.

**Proposition 20.6** *Assume that the matrix $\Sigma_T$ is invertible. The log-likelihood in (20.7) is maximized by taking*

*(i)  $W = [\lambda_1 e_1, \ldots, \lambda_p e_p]$ where $e_1, \ldots, e_p$ are the eigenvectors of $\Sigma_T$ associated to the $p$ largest eigenvalues, and $\lambda_i = \sqrt{\delta_i^2 - \sigma^2}$, where $\delta_i^2$ is the eigenvalue of $\Sigma$ associated to $e_i$;*

*(ii)  and*

$$\sigma^2 = \frac{1}{d-p}\sum_{i=p+1}^{d}\delta_i^2.$$

PROOF  We make the following change of variables: let $\rho^2 = 1/\sigma^2$ and

$$\mu_i^2 = \frac{1}{\sigma^2} - \frac{1}{\lambda_i^2 + \sigma^2}.$$

Let $Q = [\mu_1 e_1, \ldots, \mu_p e_p]$. We have

$$(WW^T + \sigma^2\mathrm{Id})^{-1} = \rho^2\mathrm{Id} - QQ^T.$$

To see this, complete $(e_1, \ldots, e_p)$ into an orthonormal basis of $\mathbb{R}^d$, letting $e_{p+1}, \ldots, e_d$ denote the added vectors. Then

$$WW^T + \sigma^2\mathrm{Id} = \sum_{i=1}^{p}(\lambda_i^2 + \sigma^2)e_i e_i^T + \sum_{i=p+1}^{d}\sigma^2 e_i e_i^T$$

so that

$$(WW^T + \sigma^2\mathrm{Id})^{-1} = \sum_{i=1}^{p}(\lambda_i^2 + \sigma^2)^{-1}e_i e_i^T + \sum_{i=p+1}^{d}\sigma^{-2}e_i e_i^T = \rho^2\mathrm{Id} - QQ^T.$$

Using these variables, we can reformulate the problem as the minimization of

$$-\sum_{i=1}^{p}\log(\rho^2 - \mu_i^2) - (d-p)\log\rho^2 + \rho^2\mathrm{trace}(\Sigma) - \sum_{j=1}^{p}\mu_j^2 e_j^T \Sigma e_j.$$

From theorem 2.3, we have

$$\sum_{j=1}^{p} \mu_j^2 e_j^T \Sigma e_j \le \sum_{j=1}^{p} \mu_j^2 \delta_j^2$$

and this upper bound is attained by letting $e_1, \ldots, e_p$ be the first $p$ eigenvectors of $\Sigma$. Using this, we see that $\sigma^2, \mu_1^2, \ldots, \mu_p^2$ must minimize

$$-\sum_{i=1}^{p} \log(\rho^2 - \mu_i^2) - (d-p)\log \rho^2 + \rho^2 \sum_{j=1}^{d} \delta_j^2 - \sum_{j=1}^{p} \mu_j^2 \delta_j^2.$$

Computing the solution is elementary and left to the reader, and yields, when expressed as functions of $\sigma^2, \lambda_1^2, \ldots, \lambda_p^2$, the expressions given in the statement of the theorem. ∎

## 20.4 Generalized PCA

We now discuss a dimension reduction method called generalized PCA (GPCA) [200] that, instead of looking for the best linear approximation of the training set by one specific subspace, provides an approximation by a finite union of such spaces.

As a motivation, consider the situation in fig. 20.1 in which part of the data is aligned along one direction in space, and another part along another direction. Then, the only information that PCA can retrieve (provided that the two directions intersect) is the plane generated by the two directions, which will be captured by the two principal components. PCA will not be able to determine the individual directions. GPCA addresses this type of situation as follows.



Figure 20.1: PCA cannot distinguish between the situations depicted in the two datasets.

For simplicity, assume that we are trying to decompose the data along unions of hyperplanes in $\mathbb{R}^d$. Such hyperplanes have equations of the form $u^T \tilde{x} = 0$ where $\tilde{x}$ is our notation for the vector $(1, x^T)^T$. If we have two hyperplanes, specified by $u_1$ and $u_2$ and all the training samples approximately belong to one of them, then one has, for all $k = 1, \ldots, N$:

$$(u_1^T \tilde{x}_k)(u_2^T \tilde{x}_k) = \tilde{x}_k^T u_1 u_2^T \tilde{x}_k \simeq 0.$$

Similarly, for $n$ hyperplanes, the identity is, for $k = 1, \ldots, N$:

$$\prod_{j=1}^{n} (u_j^T \tilde{x}_k) \simeq 0.$$

Write

$$\prod_{j=1}^{n} (u_j^T x) = \sum_{1 \leq i_1, \ldots, i_n \leq d} u_1(i_1) \cdots u_n(i_n) x^{(i_1)} \cdots x^{(i_n)}$$

in the form (by regrouping the terms associated with the same powers of $x$)

$$F(x) = \sum_{p_1 + \ldots + p_d = n} q_{p_1 \ldots p_d} (x^{(1)})^{p_1} \ldots (x^{(d)})^{p_d}. \tag{20.8}$$

The collection of $\binom{n+d-1}{n}$ numbers $Q = (q_{p_1 \ldots p_n}, p_1 + \cdots + p_d = n)$ takes a specific form (that we will not need to make explicit) as a function of the unknown $u_1, \ldots, u_n$, but the first step of GPCA ignores this constraint and estimates $Q$ minimizing

$$\sum_{k=1}^{N} \left( \sum_{p_1 + \ldots + p_d = n} q_{p_1 \ldots p_d} (x_k^{(1)})^{p_1} \ldots (x_k^{(d)})^{p_d} \right)^2$$

under the constraint $\sum q_{p_1 \ldots p_n}^2 = 1$ (to avoid trivial solutions). Choosing an ordering on the set of indices $(p_1, \ldots, p_d)$ such that $p_1 + \cdots p_d = n$, one can stack the coefficients in $Q$ and the monomials $(x_k^{(1)})^{p_1} \ldots (x_k^{(d)})^{p_d}$ to form two vectors denoted $Q$ (with some abuse of notation) and $V(x_k)$. One can then rewrite the problem of determining $Q$ as minimizing $Q^T \Sigma Q$ subject to $|Q|^2 = 1$, where

$$\Sigma = \sum_{k=1}^{N} V(x_k) V(x_k)^T.$$

The solution is given by the eigenvector associated with the smallest eigenvalue of $\Sigma$. If the model is exact, this eigenvalue should be zero, and if only one decomposition of the data in a set of distinct hyperplanes exists (i.e., if $n$ is not chosen too large), then $Q$ is the unique solution up to a multiplicative constant.

Once $Q$ is found, it remains to identify the vectors $u_1, \ldots, u_n$. This identification can be obtained by inspecting the gradient of $F$ on the union of hyperplanes. Indeed, one has, for $x \in \mathbb{R}^d$,

$$\nabla F(x) = \sum_{j=1}^{n} \left( \prod_{j' \neq j} u_{j'}^T x \right) u_j$$

However, if $x$ belong in one and only one of the hyperplanes, say $x^T u_j = 0$, then all terms in the sum vanish but one and $\nabla F(x)$ is proportional to $u_j$. So, if the model is exact, one has, for each $k = 1, \ldots, N$, either $\nabla F(x_k) = 0$ (if $x_k$ belongs to the intersection of two hyperplanes) or $\nabla F(x_k)/|\nabla F(x_k)| = \pm u_j$ for some $j$, and the sign ambiguity can be removed by ensuring, for example, that the first non-vanishing coordinate of $u_j$ is positive. (The gradient of $F$ can be computed from $Q$ using (20.8).) The computation of $\nabla F$ on training data therefore allows for an exact computation of the hyperplanes.

In practice, when noise is present, one cannot expect this computation to be exact. The vectors $u_1, \ldots, u_n$ can be estimated by clustering the collection of non-vanishing gradients $\nabla F(x_k)$, $k = 1, \ldots, N$. For example, one can compute a dissimilarity matrix such as $d_{kl} = 1 - \cos^2(\theta_{kl})$, where $\theta_{kl}$ is the angle between $\nabla F(x_k)$ and $\nabla F(x_l)$, and apply one of the methds discussed in section 19.4.1.

This analysis provides a decomposition of the training set into $n$ (or fewer) hyperplanes. The computation can then be recursively refined in order to obtain smaller dimensional subspaces by applying the same method separately to each hyperplane.

## 20.5 Nuclear norm minimization and robust PCA

### 20.5.1 Low-rank approximation

One can also interpret PCA in terms of low-rank matrix approximations. Let $\mathcal{X}_c$ be the $N$ by $d$ matrix $(x_1 - \bar{x}, \ldots, x_N - \bar{x})^T$, which, in generic situations, has rank $d - 1$. Then PCA with $p$ components is equivalent to minimizing, over all $N$ by $d$ matrices $\mathcal{Z}$ of rank $p$, the norm of the difference

$$|\mathcal{X}_c - \mathcal{Z}|^2 = \text{trace}((\mathcal{X}_c - \mathcal{Z})^T (\mathcal{X}_c - \mathcal{Z})). \tag{20.9}$$

The quantity $|A|^2 = \text{trace}(A^T A)$ is the sum of square of the entries of $A$, which is often referred to as the (squared) Frobenius norm. We have

$$|A|^2 = \sum_{k=1}^{d} \sigma_k^2$$

where $\sigma_1, \ldots, \sigma_d$ are the singular values of $A$, i.e., the square roots of the eigenvalues of $A^T A$.

We first note the following characterization of rank-$p$ matrices.

**Proposition 20.7** *A matrix $\mathcal{Z}$ has rank $p$ if and only if it can be written in the form $\mathcal{Z} = AW^T$ where $A$ is $N \times p$, and $W$ is $d \times p$ with $W^T W = \mathrm{Id}_{\mathbb{R}^p}$, i.e., $W = [e_1, \ldots, e_p]$ where the columns form an orthonormal family of $\mathbb{R}^d$.*

PROOF The "if" part is obvious and we prove the "only if" part. Assume that $\mathcal{Z}$ has rank $p$. Take $W = [e_1, \ldots, e_p]$, where $(e_1, \ldots, e_p)$ is an orthonormal family in $\mathrm{Null}(\mathcal{Z})^\perp$. Letting $e_{p+1}, \ldots, e_d$ denote an orthonormal basis of $\mathrm{Null}(\mathcal{Z})$, we have $\sum_{i=1}^{d} e_i e_i^T = \mathrm{Id}_{\mathbb{R}^d}$ and

$$\mathcal{Z} = \mathcal{Z} \sum_{i=1}^{d} e_i e_i^T = \mathcal{Z} \sum_{i=1}^{p} e_i e_i^T = \mathcal{Z} W W^T$$

so that one can take $A = \mathcal{Z}W$.                                     ■

Using this representation and letting $z_k^T$ be the $k$th row vector of $\mathcal{Z}$, we have

$$|\mathcal{X}_x - \mathcal{Z}|^2 = \sum_{k=1}^{N} |x_k - \bar{x} - z_k|^2 = \sum_{k=1}^{N} \left| x_k - \bar{x} - \sum_{j=1}^{p} a_k^{(j)} e_j \right|^2.$$

With fixed $e_1, \ldots, e_p$, the optimal matrix $A$ has coefficients $a_k^{(j)} = (x_k - \bar{x})^T e_j$. In matrix form, this is:

$$\mathcal{Z} = \mathcal{X}_c \left( \sum_{j=1}^{p} e_j e_j^T \right).$$

We therefore retrieve the PCA formulation that we gave in section 20.1, in the special case of $H = \mathbb{R}^d$ with the standard Euclidean product. The lowest value achieved by the PCA solution is

$$|\mathcal{X}_c - \mathcal{Z}|^2 = N \sum_{k=p+1}^{d} \lambda_k^2$$

where $\lambda_1^2, \ldots, \lambda_d^2$ are the eigenvalues of the covariance matrix computed from $x_1, \ldots, x_N$, who are also the squared singular values of the matrix $\mathcal{X}_c$ divided by $N$.

In this section, we will explore variations on PCA in which the minimization of $|\mathcal{X}_c - \mathcal{Z}|^2$ is completed with a penalty that depends on the singular values of the matrix $\mathcal{Z}$. As a first example, one can modify PCA by adding a penalty on the rank (i.e., on the number of non-zero singular values), minimizing:

$$\gamma |\mathcal{X}_c - \mathcal{Z}|^2 + \mathrm{rank}(\mathcal{Z})$$

for some parameter $\gamma > 0$. However, the solution to this problem is a small variation of that of standard PCA. It is indeed given by standard PCA with $p$ components where $p$ minimizes

$$N\gamma \sum_{k=p+1}^{d} \lambda_k^2 + p = N\gamma \sum_{k=p+1}^{d} (\lambda_k^2 - (N\gamma)^{-1}) + d,$$

i.e., $p$ is the index of the last eigenvalue that is larger than $(N\gamma)^{-1}$.

### 20.5.2 The nuclear norm

Based on the fact that $\mathrm{rank}(\mathcal{Z})$ is the number of non-zero singular values of $\mathcal{Z}$, one can use the same heuristic as in the development of the lasso, and replace counting the non-zero values by the sum of the absolute values of the singular values, which is just the sum of singular values since they are non-negative. This provides the nuclear norm of $A$, defined in section 2.4 by

$$|A|_* = \sum_{k=1}^{d} \sigma_k$$

where $\sigma_1, \ldots, \sigma_d$ are the singular values of $A$. We will consider below the problem of minimizing

$$\gamma|\mathcal{X}_c - \mathcal{Z}|^2 + |\mathcal{Z}|_* \tag{20.10}$$

and show that its solution is once again similar to PCA.

We recall the characterization of the nuclear norm proposition 2.6. If $A$ is an $N$ by $d$ matrix,

$$|A|_* = \max\left\{ \mathrm{trace}(UAV^T) : U \text{ is } N \times N \text{ and } U^T U = \mathrm{Id}, V \text{ is } d \times d \text{ and } V^T V = \mathrm{Id} \right\}.$$

In Cai et al. [45], the authors consider the minimization of (20.10) and prove the following result. Recall that we have defined the shrinkage function $S_\tau : t \mapsto \mathrm{sign}(t)\max(|t| - \tau, 0)$ (with $\tau \geq 0$), using the same notation $S_\tau(X)$ when applying $S_\tau$ to every entry of a vector or matrix $X$. Following Cai et al. [45], we define the *singular value thresholding operator* $A \mapsto \mathbb{S}_\tau(A)$, where $A$ is any rectangular matrix, by

$$\mathbb{S}_\tau(A) = U S_\tau(\Delta) V^T$$

when $A = U\Delta V^T$ is a singular value decomposition of $A$.

**Proposition 20.8** *Let us assume without loss of generality that $N \geq d$. The function $\mathcal{Z} \mapsto \gamma|\mathcal{X}_c - \mathcal{Z}|^2 + |\mathcal{Z}|_*$ is minimized by $\mathcal{Z} = \mathbb{S}_{1/2\gamma}(\mathcal{X})$.*

PROOF Representing $\mathcal{Z}$ by its singular value decomposition, we have the equivalent formulation of minimizing

$$
\begin{aligned}
F(U,V,D) &= \gamma|\mathcal{X}_c - UDV^T|^2 + |D|_* \\
&= \gamma|\mathcal{X}_c|^2 - 2\gamma\,\mathrm{trace}(\mathcal{X}_c^T UDV^T) + \gamma|D|^2 + |D|_*
\end{aligned}
$$

over all orthonormal matrices $U$ and $V$ and diagonal matrices with non-negative coefficients $D$. From theorem 2.1, we know that $\mathrm{trace}(\mathcal{X}_c^T UDV^T)$ is less than the sum of the products of the non-increasingly ordered singular values of $\mathcal{X}_c$ and $D$ and this upper bound is attained by taking $U = \bar{U}$ and $V = \bar{V}$ where $\bar{U}$ and $\bar{V}$ are the matrices providing the SVD of $\mathcal{X}_c$, i.e., such that $\mathcal{X}_c = \bar{U}\Delta\bar{V}^T$ where $\Delta$ is diagonal with non-decreasing coefficients along the diagonal. So, letting $\lambda_1 \geq \cdots \geq \lambda_d \geq 0$ and $\mu_1 \geq \cdots \geq \mu_d \geq 0$ be the singular values of $\mathcal{X}_c$ and $\mathcal{Z}$, we have just proved that, for any $D$,

$$
F(U,V,D) \geq F(\bar{U},\bar{V},D) = -2\gamma\sum_{i=1}^{d}\mu_i\lambda_i + \gamma\sum_{i=1}^{d}\mu_i^2 + \sum_{i=1}^{d}\mu_i.
$$

The lower bound is minimized when $\mu_i = \max(\lambda_i - 1/2\gamma, 0)$. This proves the proposition. ∎

### 20.5.3   Robust PCA

As a consequence, the nuclear norm penalty provides the same principal directions (after replacing $\gamma$ by $2\gamma$) as the rank penalty, but applies a shrinking operation rather than thresholding on the singular values. The difference is however more fundamental if, in addition to using the nuclear norm as a penalty, on replaces the squared Frobenius norm on the approximation error by the $\ell^1$ norm, where, for an $n$ by $m$ matrix $A$ with coefficients $(a(i,j))$,

$$
|A|_{\ell_1} = \sum_{i,j}|a(i,j)|.
$$

This is the formulation of robust PCA [49], which minimizes

$$
\gamma|\mathcal{X}_c - \mathcal{Z}|_{\ell^1} + |\mathcal{Z}|_* \tag{20.11}
$$

with respect to $\mathcal{Z}$.

Robust PCA (which was initially named Principal Component Pursuit by the authors in Candès et al. [49]) is designed for situations in which $\mathcal{X}_c$ can be decomposed as the sum of a low-rank matrix $\mathcal{Z}$ and of a sparse residual $S$. Some theoretical justification was provided in the original paper, stating that if $\mathcal{X}_c = \mathcal{Z}+S$, with $\mathcal{Z} = UDV^T$ (its singular value decomposition) such that $U$ and $V$ are sufficiently "diffuse" and

rank($\mathcal{Z}$) is small enough, with the residual's sparsity pattern taken uniformly at random over the subsets of entries of $S$ with a sufficiently small cardinality, then robust PCA is able to reconstruct the decomposition exactly with high probability (relative to the random selection of the sparsity pattern of $S$). We refer to Candès et al. [49] for the long proof that justifies this statement.

Robust PCA can be solved using the ADMM algorithm (section 3.5.5) after reformulating the problem as the minimization of

$$\gamma |R|_{\ell_1} + |\mathcal{Z}|_*$$

subject to $R + \mathcal{Z} = \mathcal{X}_c$. The algorithm therefore iterates over the following steps.

$$\begin{cases} \mathcal{Z}^{(k+1)} = \underset{\mathcal{Z}}{\operatorname{argmin}} \left( |\mathcal{Z}|_* + \frac{1}{2\alpha}|\mathcal{Z} + R^{(k)} - \mathcal{X}_x + U^{(k)}|^2 \right) \\ R^{(k+1)} = \underset{R}{\operatorname{argmin}} \left( \gamma|R|_{\ell^1} + \frac{1}{2\alpha}|\mathcal{Z}^{(k+1)} + R - \mathcal{X}_c + U^{(k)}|^2 \right) \\ U^{(k+1)} = U^{(k)} + \mathcal{Z}^{(k+1)} + R^{(k+1)} - \mathcal{X}_c \end{cases} \qquad (20.12)$$

The first minimization is covered by proposition 2.6 and yields

$$\mathcal{Z}^{(k+1)} = \mathbb{S}_\alpha(\mathcal{X}_c - R^{(k)} - U^{(k)}).$$

The second minimization is solved by a standard shrinking operation, i.e.,

$$R^{(k+1)} = S_{\gamma\alpha}(\mathcal{X}_c - \mathcal{Z}^{(k+1)} - U^{(k)}).$$

Using this, we can rewrite the robust PCA algorithm as the sequence of fairly simple iterations.

---

**Algorithm 20.1**

(1) Choose a small enough constant $\alpha$ and a very small tolerance level $\epsilon$.

(2) Initialize the algorithm with $N$ by $d$ matrices $R^{(0)}$ and $U^{(0)}$ (e.g., equal to zero).

(3) At step $n$, apply the iteration:

$$\begin{cases} \mathcal{Z}^{(k+1)} = \mathbb{S}_\alpha(\mathcal{X}_c - R^{(k)} - U^{(k)}) \\ R^{(k+1)} = S_{\gamma\alpha}(\mathcal{X}_c - \mathcal{Z}^{(k+1)} - U^{(k)}) \\ U^{(k+1)} = U^{(k)} + \mathcal{Z}^{(k+1)} + R^{(k+1)} - \mathcal{X}_c \end{cases} \qquad (20.13)$$

(4) Stop the algorithm is the variation compared to variables at the previous step is below the tolerance level. Otherwise, apply step $n + 1$.

---

## 20.6   Independent component analysis

Independent component analysis (ICA) is a factor analysis method that represents a $d$-dimensional random variable $X$ in the form $X = AY$ where $A$ is a fixed $d \times d$ invertible matrix and $Y$ is a $d$-dimensional random vector with independent components. There are two main approaches in this setting. The first one optimizes the matrix $W = A^{-1}$ so that the components of $WX$ are "as independent as possible" according to a suitable criterion. The second one is model-based, where a statistical model is assumed for $Y$, and its parameters, together with the entries of the matrix $A$, are estimated via maximum likelihood. Before describing each of these methods, we first discuss the extent to which the coefficients of $A$ are identifiable.

### 20.6.1   Identifiability

A statistical model is identifiable if its parameters (which could be finite- of infinite-dimensional) are uniquely defined by the distribution of the observable variables. In the case of ICA, this question boils down to deciding whether $AY \sim A'Y'$ (i.e., they have the same probability distribution) implies that $A = A'$ (where $Y$ and $Y'$ are two random vectors with independent components).

It should be clear that the answer to this question is negative, because there are trivial transformations of the matrix $A$ that do not break the ICA model. One can, for example, take any invertible diagonal matrix, $D$, and let $A' = AD^{-1}$ and $Y' = DY$. The same statement can be made if $D$ is replaced by a permutation matrix, $P$, which reorders the components of $Y$. So we know that $AY \sim A'Y'$ is possible already when $A' = ADP$ where $D$ is diagonal and invertible and $P$ is a permutation matrix. Note that iterating such matrices (i.e., letting $A' = ADPD'P'$) does not extend the class of transformations because one has $DP = PP^{-1}DP$ and one can easily check that $P^{-1}DP$ is diagonal, so that one can rewrite any product of permutations and diagonal matrices as a single diagonal matrix multiplied by a single permutation.

It is interesting, and fundamental for the well-posedness of ICA, that, under one important additional assumption, the indeterminacy in the identification of $A$ stops at these transformations. The additional assumption is that at most one of the components of $Y$ follows a Gaussian distribution. That such a restriction is needed is clear from the fact that one can transform any Gaussian vector $Y$ with independent components into another, $BY$, one as soon as $BB^T$ is diagonal. If two or more components of $Y$ are Gaussian, one can restrict these matrices $B$ to only affect those components. If only one of them is Gaussian, such an operation has no effect.

The following theorem is formally stated in Comon [54], and is a rephrasing of the Darmois-Skitovitch theorem [57, 179]. The proof of this theorem relies on complex analysis arguments on characteristic functions and is beyond the scope of

these notes (see Kagan et al. [101] for more details).

**Theorem 20.9** *Assume that $Y$ is a random vector with independent components, such that at most one of its components is Gaussian. Let $A$ be an invertible linear transformation and $\tilde{Y} = CY$. Then the following statements are independent.*

*(i) For all $i \neq j$, the components $\tilde{Y}^{(i)}, \tilde{Y}^{(j)}$ are independent.*

*(ii) $\tilde{Y}^{(1)}, \ldots, \tilde{Y}^{(d)}$ are mutually independent.*

*(iii) $C = DP$ is the product on a diagonal matrix and of a permutation.*

The equivalence of (ii) and (iii) implies that the ICA model is identifiable up to multiplication on the right by a permutation and a diagonal matrix. Indeed, if $X = AY = A'Y'$ are two decompositions, then it suffices to apply the theorem to $C = (A')^{-1}A$ to conclude. The equivalence of (i) and (ii) is striking, and has the important consequence that, if the data satisfies the ICA model, then, in order to identify $A$ (up to the listed indeterminacy), it suffices to look for $Y = A^{-1}X$ with pairwise independent components, which is a much lesser constraint than full mutual independence.

As a final remark on the Gaussian indeterminacy, we point out that, if the mean ($m$) and covariance matrix ($\Sigma$) of $X$ are known (or estimated from data), the ICA problem can be reduced to looking for orthogonal transformations $A$. Indeed, assuming $X = AY$ and letting $\tilde{X} = \Sigma^{-1/2}(X - m)$ and $\tilde{Y} = D^{-1/2}(Y - A^{-1}m)$, where $D$ is the (diagonal) covariance matrix of $Y$, we have

$$\tilde{X} = \Sigma^{-1/2}(AY - m) = \Sigma^{-1/2}AD^{1/2}\tilde{Y}.$$

Letting $\tilde{A} = \Sigma^{-1/2}AD^{1/2}$, we have $\mathrm{Id}_{\mathbb{R}^d} = E(\tilde{X}\tilde{X}^T) = \tilde{A}\tilde{A}^T$ so that $\tilde{A}$ is orthogonal. This shows that the ICA problem for $\tilde{X}$ in the form $\tilde{X} = \tilde{A}\tilde{Y}$ with the restriction that $\tilde{A}$ is orthogonal has a solution, and also provides a solution of the original ICA problem by letting $A = \Sigma^{1/2}\tilde{A}$ and $Y = \tilde{Y} - \tilde{A}^{-1}\Sigma^{-1/2}m$. Therefore, the indeterminacy associated with Gaussian vectors is as general as possible up to a normalization of first and second moments.

### 20.6.2 Measuring independence and non-Gaussianity

Independence between $d$ variables is a very strong property and its complete characterization is computationally challenging. The fact that the joint p.d.f.of the $d$ variables (we will restrict, to simplify our discussion, to variables that are absolutely continuous) factorizes into the product of the marginal p.d.f.'s of each variable can be measured by computing the mutual information between the variables, defined

by (letting $\varphi_Z$ denote the p.d.f. of a variable $Z$)

$$I(Y) = \int \frac{\varphi_Y(y)}{\prod_{i=1}^{d} \varphi_{Y^{(i)}}(y^{(i)})} \varphi_Y(y) dy.$$

The mutual information is always non-negative and vanishes only if the components of $Y$ are mutually independent. Therefore, one can represent ICA as an optimization problem minimizing $I(WX)$ with respect to all invertible matrices $W$ (so that $W = A^{-1}$). Letting

$$h(Y) = -\int \log \varphi_Y(y) \varphi_Y(y) dy$$

denote the "differential entropy" of $Y$, we can write

$$I(Y) = \sum_{i=1}^{d} h(Y^{(i)}) - h(Y).$$

If $Z = WX$, then $\varphi_Z(z) = \varphi_X(W^{-1}x)|\det(W)|^{-1}$. Using this expression in $h(Z)$ and making a change of variables yields $h(WZ) = h(X) + \log|\det W|$ and

$$I(WX) = \sum_{i=1}^{d} h(Z^{(i)}) - \log|\det(W)| - h(X).$$

This shows that the optimal $W$ can be obtained by minimizing

$$F(W) = \sum_{i=1}^{d} h(W^{(i)}X) - \log|\det(W)|$$

where $W^{(i)}$ is the $i$th row of $W$. This brings a notable simplification, since this expression only involves differential entropies of scalar variables, but still remains a challenging problem.

In Comon [54], it is proposed to use cumulant expansions of the entropy around that of a Gaussian with identical mean and variance to approximate the differential entropy. If $\xi \sim \mathcal{N}(m, \sigma^2)$, then

$$h(\xi) = \frac{1}{2} + \frac{1}{2} \log(2\pi\sigma^2).$$

Define, for a general random variable $U$ with standard deviation $\sigma_U$, the non-Gaussian entropy, or negentropy, defined by

$$v(U) = \frac{1}{2} + \frac{1}{2} \log(2\pi\sigma_U^2) - h(U).$$

One can shows that $v(U) \geq 0$ and is equal to 0 if and only if $U$ is Gaussian. One can rewrite $F(W)$ as

$$F(W) = \frac{d}{2} + \frac{d}{2}\log(2\pi) + \sum_{i=1}^{d}\log(\sigma^2_{W^{(i)}X}) - \sum_{i=1}^{d} v(W^{(i)}X) - \log|\det(W)|$$

As we remarked earlier, if we replace $X$ by $\Sigma^{-1/2}(X - m)$ (after estimating the covariance matrix of $X$), there is no loss of generality in requiring that $W$ is an orthogonal matrix, in which case both $\sigma^2_{W^{(i)}X}$ and $|\det W|$ are equal to 1. Assuming such a reduction is done, we see that the problem now requires to maximize

$$\sum_{i=1}^{d} v(W^{(i)}X) \tag{20.14}$$

among all orthogonal matrices $W$. Still in Comon [54], an approximation of the negentropy $v(U)$ is provided as a function of the third and fourth cumulants of the distribution of $U$. These are given by

$$\kappa_3 = E((U - E(U))^3)$$

and

$$\kappa_4 = E((U - E(U))^4) - 3\sigma^4_U.$$

In particular, when $U$ is normalized, i.e., $E(U) = 0$ and $\sigma^2_U = 1$, we have $\kappa_3 = E(U^3)$ and $\kappa_4 = E(U^4) - 3$. Under the same assumption, it is proposed in Comon [54] to use the approximation

$$v(U) \sim \frac{\kappa_3^2}{12} + \frac{\kappa_4^2}{48} + \frac{7\kappa_3^4}{48} - \frac{\kappa_3^2\kappa_4}{8}.$$

This approximation was derived from an Edgeworth expansion of the p.d.f. of $U$, which can be seen as a Taylor expansion around a Gaussian distribution. Plugging this expression into (20.14) provides an expression that can be maximized in $W$ where the cumulants are replaced by their sample estimates. However, the maximized function involves high-degree polynomials in the unknown coefficients of $W$, and this simplified problem still presents numerical challenges.

An alternative approximation of the negentropy has been proposed in Hyvärinen [94] relying on the maximum entropy principle, described in the following theorem. Associate to any random variable $Y : \mathcal{G} \to \mathbb{R}$ the differential entropy

$$h_\mu(Y) = -\int_{\mathcal{G}} \log\varphi_Y(x)\varphi_Y(x)d\mu(x)$$

if the distribution of $Y$ has a density, denoted $\varphi_Y$, with respect to $\mu$ and $h_\mu(Y) = -\infty$ otherwise. Use also the same notation

$$h_\mu(\varphi) = -\int_{\mathcal{G}} \log \varphi(x)\varphi(x)d\mu(x)$$

for a p.d.f. $\varphi$ with respect to $\mu$ (i.e., such that $\varphi$ is non-negative and has integral 1). Then, the following is true.

**Theorem 20.10** *Let $g = (g^{(1)},\ldots,g^{(p)})^T$ be a function defined on a measurable space $\mathcal{G}$, taking values in $\mathbb{R}^p$, and let $\mu$ be a measure on $\mathcal{G}$. Let $\Gamma_\mu$ be the set of all $\lambda = (\lambda^{(1)},\ldots,\lambda^{(p)}) \in \mathbb{R}^p$ such that*

$$\int_{\mathcal{G}} \exp\left(\lambda^T g(y)\right) d\mu(y) < \infty. \tag{20.15}$$

*Then*

$$h_\mu(Y) \le \inf\left\{-\lambda^T E(g(Y)) + \log \int_{\mathcal{G}} \exp\left(\lambda^T g(y)\right) d\mu(y) : \lambda \in \Gamma_\mu\right\}. \tag{20.16}$$

*Define, for $\lambda \in \Gamma_\mu$,*

$$\psi_\lambda(x) = \frac{\exp\left(\lambda^T g(x)\right) d\mu(x)}{\int_{\mathcal{G}} \exp\left(\lambda^T g(y)\right) d\mu(y)}. \tag{20.17}$$

*Assume that the infimum in (20.16) is attained at an interior point $\lambda^*$ of $\Gamma_\mu$. Then*

$$h(\varphi_{\lambda^*}) = \max\{h(\tilde{Y}) : E_{\tilde{Y}}(g) = E_Y(g), i = 1,\ldots,p\}. \tag{20.18}$$

PROOF Let $Y$ be a random variable with p.d.f. $\varphi_Y$ with respect to $\mu$ (otherwise the lower bound in (20.16) is $-\infty$). Then

$$h_\mu(Y) + \lambda E(g(Y)) - \log \int \exp\left(\lambda g(y)\right) d\mu(y) = -\int_{\mathcal{G}} \varphi_Y(x) \log \frac{\varphi_Y(x)}{\psi_\lambda(x)} d\mu(x) \le 0$$

since $\int_{\mathcal{G}} \varphi_Y(x) \log \frac{\varphi_Y(x)}{\psi_\lambda(x)} d\mu(x)$ is a KL divergence and is always non-negative.

Assume that $\lambda$ is in $\mathring{\Gamma}_\mu$. Then, there exists $\epsilon > 0$ such that, for any $u \in \mathbb{R}^p$, $|u| = 1$, $\lambda + \epsilon u \in \Gamma_\mu$. Using the fact that $e^\beta \ge e^\alpha + (\beta - \alpha)e^\alpha$, we can write

$$\epsilon u^T g e^{\lambda^T g} \le e^{(\lambda+\epsilon u)^T g} - e^{\lambda^T g}$$
$$-\epsilon u^T g e^{\lambda^T g} \le e^{(\lambda-\epsilon u)^T g} - e^{\lambda^T g}$$

yielding

$$\epsilon |u^T g| e^{\lambda^T g} \le \max(e^{(\lambda+\epsilon u)^T g}, e^{(\lambda-\epsilon u)^T g}) - e^{\lambda^T g} \le e^{(\lambda+\epsilon u)^T g} + e^{(\lambda-\epsilon u)^T g} - e^{\lambda^T g}.$$

Since the upper-bound is integrable with respect to $\mu$, so is the lower bound, showing that (taking $\boldsymbol{u}$ in the canonical basis of $\mathbb{R}^p$)

$$\int_{\mathcal{G}} |g^{(i)}(y)| e^{\lambda^T g(y)} d\mu(y) < \infty$$

for all $i$, or

$$\int_{\mathcal{G}} |\boldsymbol{g}(y)| e^{\lambda^T g(y)} d\mu(y) < \infty.$$

Let $\boldsymbol{c} = E(\boldsymbol{g}(Y))$ and define

$$\Psi_c(\lambda) = -\boldsymbol{c}^T \lambda + \log \int \exp(\lambda^T \boldsymbol{g}(y)) dy. \tag{20.19}$$

Then

$$\partial_\lambda \Psi_c = -c^T + \frac{\int \boldsymbol{g}(x)^T \exp(\lambda^T \boldsymbol{g}(x)) dx}{\int \exp(\lambda^T \boldsymbol{g}(y)) dy} = -c^T + \int_{\mathcal{G}} \boldsymbol{g}(x)^T \psi_\lambda(x) dx.$$

Since $\lambda^*$ is a minimizer, we find that, if $\tilde{Y}$ is a random variable with p.d.f. $\lambda^*$, then $E(\tilde{Y}) = c = E(Y)$. In that case, the upper-bound in (20.16) is $h_\mu(\tilde{Y})$, proving (20.18). ∎

**Remark 20.11** The previous theorem is typically applied with $\mu$ equal to Lebesgue's measure on $\mathcal{G} = \mathbb{R}^d$ or to a counting measure with $\mathcal{G}$ finite. To rewrite the statement of theorem 20.10 in those cases, it suffices to replace $d\mu(x)$ by $dx$ for the former, and integrals by sums over $\mathcal{G}$ for the latter. In the rest of the discussion, we restrict to the case when $\mu$ is Lebesgue's measure, using $h(Y)$ instead of $h_\mu(Y)$. ◆

**Remark 20.12** This principle justifies, in particular, that the negentropy is always non-negative since it implies that a distribution that maximizes the entropy given its first and second moments must be Gaussian. ◆

The right-hand side of (20.16) provides a variational approximation of the entropy. If one uses this approximation when minimizing $h(W^{(1)}X) + \cdots + h(W^{(d)}X)$, the resulting problem can be expressed as a minimization, with respect to $W$ and $\lambda_1, \ldots, \lambda_d \in \mathbb{R}^p$ of

$$-\sum_{j=1}^d \lambda^T E(\boldsymbol{g}(W^{(j)}X)) + \sum_{j=1}^d \log \int \exp\left(\lambda^T \boldsymbol{g}(y)\right) dy.$$

While it would be possible to solve this optimization problem directly, a further approximation of the upper bound can be developed leading to a simpler procedure.

We have seen in the previous proof that, defining $\Psi_c$ by (20.19) and denoting by $E_\lambda$ the expectation with respect to $\varphi_\lambda$, one has

$$\nabla \Psi_c(\lambda) = -c + E_\lambda(g).$$

Taking the second derivative, one finds

$$\nabla^2 \Psi_c(\lambda) = E_\lambda((g - E_\lambda(g))(g - E_\lambda(g))^T).$$

Now choose $c_0$ such that a maximizer of $\Psi_{c_0}(\lambda)$, say, $\lambda^{c_0}$, is known. If $c$ is close to $c_0$, a first order expansion indicates that, for $\lambda^c$ maximizing $\Psi_c$, one should have

$$\lambda^c \simeq \lambda^{c_0} + \nabla^2 \Psi_c(\lambda^{c_0})^{-1}(c - c_0)$$

with

$$\Psi_c(\lambda^c) \simeq \Psi_c(\lambda^{c_0}) - (c - c_0)^T \nabla^2 \Psi_c(\lambda^{c_0})^{-1}(c - c_0).$$

One can then use the right-hand side as an approximation of the optimal entropy.

This leads to simple computations under the following assumptions. First, assume that the first two functions $g^{(1)}$ and $g^{(2)}$ are $u$ and $u^2/\sqrt{3}$. Let $\varphi_0$ be the p.d.f. of a standard Gaussian. Assume that the functions $g^{(j)}$ are chosen so that

$$\int g^{(i)}(u)g^{(j)}(u)\varphi_0(y)dy = \delta_{ij}$$

for $i, j = 1, \ldots, p$ and such that $\int g^{(i)}(u)\varphi_0(y)dy = 0$ for $i \ne 2$. Take

$$c_0 = \int g\varphi_0(u)du$$

so that $c_0^{(1)} = 0$, $c_0^{(2)} = 1/\sqrt{3}$ and $c_0^{(i)} = 0$ for $i \ge 2$.

Then $\lambda^{c_0}$ provides, by construction, the distribution $\varphi_0$ and for any $c$, $\nabla^2 \Psi_c(\lambda^{c_0}) = \mathrm{Id}_{\mathbb{R}^p}$. With these assumptions, the approximation is

$$\Psi_c(\lambda) = h(\varphi_0) - |c - c_0|^2$$
$$= \frac{1}{2}(1 + \log 2\pi) - \sum_{j \ge 3}(c^{(j)})^2$$

(assuming that the data is centered and normalized so that $c^{(1)} = 0$ and $c^{(2)} = 1/\sqrt{3}$). The ICA problem can then be solved by maximizing

$$\sum_{j=1}^{d}\sum_{i=1}^{p} E(g^{(i)}(W^{(j)}X))^2 \qquad (20.20)$$

over orthogonal matrices $W$.

**Remark 20.13** Without the assumption made on the functions $g^{(j)}$, one needs to compute $S = \text{Cov}(g(U))^{-1}$ where $U \sim \mathcal{N}(0,1)$ and maximize

$$\sum_{j=1}^{d} (E(\boldsymbol{g}(W^{(j)}X)) - E(\boldsymbol{g}(U)))^T S (E(\boldsymbol{g}(W^{(j)}X)) - E(\boldsymbol{g}(U))).$$

Clearly, this expression can be reduced to (20.20) by replacing $\boldsymbol{g}$ by $S^{-1/2}(\boldsymbol{g} - E(\boldsymbol{g}(U)))$. Note also that we retrieve here a similar idea to the negentropy, maximizing a deviation to a Gaussian. ◆

### 20.6.3 Maximization over orthogonal matrices

In the previous discussion, we reached a few times a formulation of ICA which required optimizing a function $W \mapsto F(W)$ over all orthogonal matrices. We now discuss how such a problem may be implemented.

In all the examples that were considered, there would have been no loss of generality in requiring that $W$ is a rotation, i.e., $\det(W) = 1$. This is because one can change the sign of this determinant by simply changing the sign of one of the independent components, which is always possible. (In fact, the indeterminacy in $W$ is by right multiplication by the product of a permutation matrix and a diagonal matrix with $\pm 1$ entries.)

Let us assume that $F(W)$ is actually defined and differentiable over all invertible matrices, which form an open subset of the linear space $\mathcal{M}_d(\mathbb{R})$ of $d$ by $d$ matrices. Our optimization problem can therefore be considered as the minimization of $F$ with the constraint that $WW^T = \text{Id}_{\mathbb{R}^d}$.

Gradient descent derives from the analysis that a direction of descent should be a matrix $H$ such that $F(W + \epsilon H) < F(W)$ for small enough $\epsilon > 0$ and on the remark that $H = -\nabla F(W)$ provides such a direction. This analysis does not apply to the constrained optimization setting because, unless the constraints are linear, $W + \epsilon H$ will generally stop to satisfy the constraint when $\epsilon > 0$, requiring the use of more complex procedures. In our case, however, one can take advantage of the fact that orthogonal matrices form a group to replace the perturbation $W \mapsto W + \epsilon H$ by $W \mapsto W e^{\epsilon H}$ (using the matrix exponential) where $H$ is moreover required to be skew symmetric ($H + H^T = 0$), which guarantees that $e^{\epsilon H}$ is an orthogonal matrix with determinant 1. Now, using the fact that $e^{\epsilon H} = \text{Id} + \epsilon H + o(\epsilon)$, we can write

$$F(We^{\epsilon H}) = F(W) + \epsilon \text{trace}(\nabla F(W)^T W H) + o(\epsilon).$$

Let $\nabla^s F(W)$ be the skew symmetric part of $W^T \nabla F(W)$, i.e.,

$$\nabla^s F(W) = \frac{1}{2}(W^T \nabla F(W) - \nabla F(W)^T W).$$

Then, if $H$ is skew symmetric,

$$
\begin{aligned}
\mathrm{trace}(\nabla^s F(W)^T H) &= \frac{1}{2}\mathrm{trace}(\nabla F(W)^T W H) - \frac{1}{2}\mathrm{trace}(W^T \nabla F(W) H) \\
&= \frac{1}{2}\mathrm{trace}(\nabla F(W)^T W H) + \frac{1}{2}\mathrm{trace}(W^T \nabla F(W) H^T) \\
&= \mathrm{trace}(\nabla F(W)^T W H)
\end{aligned}
$$

so that

$$
F(W e^{\epsilon H}) = F(W) + \epsilon\,\mathrm{trace}(\nabla^s F(W)^T H) + o(\epsilon).
$$

This show that $H = -\nabla^s F(W)$ provides a direction of descent in the orthogonal group, in the sense that, if $\nabla^s F(W) \neq 0$,

$$
F(W e^{-\epsilon \nabla^s F(W)}) < F(W)
$$

for small enough $\epsilon > 0$. As a consequence, the algorithm

$$
W_{n+1} = W_n e^{-\epsilon_n \nabla^s F(W_n)}
$$

combined with a line search for $\epsilon_n$ implements gradient descent in the group of orthogonal matrices, and therefore converges to a local minimizer of $F$.

If one linearizes the r.h.s. as a function of $\epsilon$, one gets

$$
\begin{aligned}
W_n e^{-\epsilon_n \nabla^s F(W_n)} &= W_n + \frac{\epsilon_n}{2} W_n((W_n)^T \nabla F(W_n) - \nabla F(W_n)^T W_n) + o(\epsilon) \\
&= W_n + \frac{\epsilon_n}{2}(\nabla F(W_n) - W_n \nabla F(W_n)^T W_n) + o(\epsilon).
\end{aligned}
$$

As already argued, this linearized version cannot be used when optimizing over the orthogonal group. However, if one denotes by $\omega(A)$ the unitary part of the polar decomposition of $A$, i.e., $\omega(A) = (AA^T)^{-1/2}A$, then the algorithm

$$
W_{n+1} = \omega\left(W_n + \frac{\epsilon_n}{2}(\nabla F(W_n) - W_n \nabla F(W_n)^T W_n)\right)
$$

also provides a valid gradient descent algorithm.

### 20.6.4   Parametric ICA

We now describe a parametric version of ICA in which a model is chosen for the independent components of $Y$. The simplest version of to assume that all $Y^{(j)}$ are i.i.d. with some prescribed p.d.f., say, $\psi$. A typical example for $\psi$ is a logistic distribution with

$$
\psi(t) = \frac{2}{(e^t + e^{-t})^2}.
$$

If $y$ is a vector in $\mathbb{R}^d$, we will use, as usual, the notation $\psi(y) = (\psi(y^{(1)}), \ldots, \psi(y^{(d)}))^T$ for $\psi$ applied to each component of $y$.

The model parameter is then the matrix $A$, or preferably $W = A^{-1}$, and it may be estimated using maximum likelihood. Indeed, the p.d.f. of $X$ is

$$f_X(x) = |\det W| \prod_{j=1}^{d} \psi(W^{(j)}x)$$

where $W^{(j)}$ is the $j$th row of $W$, so that $W$ can be estimated by maximizing

$$\ell(W) = N \log|\det(W)| + \sum_{k=1}^{N} \sum_{j=1}^{d} \log \psi(W^{(j)}x_k).$$

If we denote by $\Gamma(W)$ the matrix with coefficients

$$\gamma_{ij}(W) = \sum_{k=1}^{N} x_k(i) \frac{\psi'(W^{(j)}x_k)}{\psi(W^{(j)}x_k)}$$

and use the fact that the gradient of $W \mapsto \log|\det W|$ is $W^{-T}$ (the inverse transpose of $W$), we can write

$$\nabla\ell(W) = N W^{-T} + \Gamma(W).$$

We need however the maximization to operate on sets of invertible matrices, and it is more natural to move in this set through multiplication than through addition, because the product of two invertible matrices is always invertible, but not necessarily their sum. So, similarly to the previous section, we will look for small variations in the form $W \mapsto We^{\epsilon H}$, or simply, in this case, $W \mapsto W(\mathrm{Id}_{\mathbb{R}^d} + \epsilon H)$. In both case, the first order expansion of the log-likelihood gives

$$\ell(W) + \epsilon \operatorname{trace}((N W^{-T} + \Gamma(W))^T W H)$$

which suggests taking

$$H = W^T(N W^{-T} + \Gamma(W)) = N\mathrm{Id} + W^T\Gamma(W).$$

Dividing $H$ by $N$, we obtain the following variant of gradient ascent for maximum likelihood

$$W_{n+1} = (1 + \epsilon_n)W_n + \epsilon_n W_n W_n^T \Gamma(W_n).$$

This algorithm numerically performs much better than standard gradient ascent. It moreover presents the advantage of avoiding computing the inverse of $W$ at each step.

### 20.6.5   Probabilistic ICA

Note that the algorithms that we discussed concerning ICA were all formulated in terms of the matrix $W = A^{-1}$, which "filters" the data into independent components. As a result, ICA requires as many independent components as the dimension of $X$. Moreover, because the components are typically normalized to have equal variance, there is no obvious way to perform dimension reduction using this method. Indeed, ICA is typically run after the data is preprocessed using PCA, this preprocessing step providing the reduction of dimension.

It is however possible to define a model similar to probabilistic PCA, assuming a limited number of components to which a Gaussian noise is added, in the form

$$X = \sum_{j=1}^{p} a_j Y^{(j)} + \sigma R$$

with $p < d$, $a_1, \ldots, a_p \in \mathbb{R}^d$, $Y^{(1)}, \ldots, Y^{(p)}$ independent variables as before, and $R \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^d})$. This model is identifiable (up to permutation and scalar multiplication of the components) as soon as none of the variables $Y^{(j)}$ is Gaussian.

Let us assume a parametric setting similar to that of the previous section, so that $Y^{(1)}, \ldots, Y^{(p)}$ are explicitly modeled as independent variables with p.d.f. $\psi$. Introduce the matrix $A = [a_1, \ldots, a_p]$, so that the model can also be written $X = AY + \sigma R$, where $A$ and $\sigma^2$ are unknown model parameters.

The p.d.f. of $X$ is now given by

$$f_X(x; A, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} \int_{\mathbb{R}^p} e^{-\frac{|x-Ay|^2}{2\sigma^2}} \left( \prod_{i=1}^{p} \psi(y^{(i)}) \right) dy^{(1)} \ldots dy^{(p)},$$

which is definitely not a closed form. Since we are in a situation in which the pair of random variables is imperfectly observed through $X$, using the EM algorithm (chapter 16) is an option, but it may, as we shall see below, lead to heavy computation. The basic step of the EM is, given current parameters $A_0, \sigma_0$, to maximize the conditional expectation (knowing $X$, for the current parameters) of the joint log-likelihood of $(X, Y)$ with respect to the new parameters. In this context, the joint distribution of $(X, Y)$ has density

$$f_{X,Y}(x, y; A, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{|x-Ay|^2}{2\sigma^2}} \prod_{i=1}^{p} \psi(y^{(i)})$$

so that, the conditional joint likelihood over the training set is

$$-\frac{Nd}{2}\log(2\pi\sigma^2)-\frac{1}{2\sigma^2}\sum_{k=1}^{N}E_{A_0,\sigma_0}(|x_k-AY|^2|X=x_k)-\sum_{k=1}^{N}\sum_{j=1}^{p}E_{A_0,\sigma_0^2}(\log\psi(Y^{(j)})|X=x_k).$$

Notice that the last term does not depend on $A, \sigma^2$, and that, given $A$, the optimal value of $\sigma^2$ is given by

$$\sigma^2 = \frac{1}{Nd}\sum_{k=1}^{N}E_{A_0,\sigma_0}(|x_k-AY|^2|X=x_k)$$

The minimization of

$$\sum_{k=1}^{N}E_{A_0,\sigma_0}(|x_k-AY|^2|X=x_k)$$

with respect to $A$ is a least square problem. Let $b_k^{(j)} = E_{A_0,\sigma_0}(Y^{(j)}|X=x_k)$ and $s_k(i,j) = E_{A_0,\sigma_0}(Y^{(i)}Y^{(j)}|X=x_k)$: the gradient of the previous term is

$$-2\sum_{k=1}^{N}E_{A_0,\sigma_0^2}((x_k-AY)Y^T|X=x_k) = -2\sum_{k=1}^{N}(x_k b_k^T - AS_k),$$

$b_k$ being the column vector with coefficients $b_k^{(j)}$ and $S_k$ the matrix with coefficients $s_k(i,j)$. The result therefore is

$$A = \left(\sum_{k=1}^{N}x_k b_k^T\right)\left(\sum_{k=1}^{N}S_k\right)^{-1}.$$

Unfortunately, the computation of the moments of the conditional distribution of $Y$ given $x_k$ (needed in $b_k$ and $S_k$) is a difficult task. The conditional density of $Y$ given $X = x_k$ is

$$g(y|x_k) = \psi(y)e^{-\frac{|A_0y-x|^2}{2\sigma_0^2}}/Z(A_0,\sigma_0)$$

from which moments cannot be computed analytically in general. Monte-Carlo sampling algorithms can be used however to approximate these moments, but they are computationally demanding. And they must be run at every step of the EM.

In place of the exact EM, one may use a mode approximation (section 16.3.1), which replaces the conditional likelihood of $Y$ given $X = x_k$ by a Dirac distribution

at the mode:

$$\hat{y}_{A_0,\sigma_0}(x_k) = \text{argmax}_y\left(\psi(y)e^{-\frac{|A_0y-x_k|^2}{2\sigma_0^2}}\right).$$

The maximization step then reduces to maximizing in $A, \sigma^2$

$$-\frac{Nd}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{k=1}^{N}\left|x_k - A\hat{y}_{A_0,\sigma_0}(x_k)\right|^2. \tag{20.21}$$

This therefore provides a two-step procedure.

---

**Algorithm 20.2 (Probabilistic ICA: mode approximation)**
   (1)  Initialize the algorithm with $A_0, \sigma_0$.

   (2)  At step $n$:

   (i)  For $k = 1, \ldots, N$, maximize $\prod_{i=1}^{p}\psi(y^{(i)})e^{-\frac{|A_ny-x_k|^2}{2\sigma_n^2}}$ to obtain $\hat{y}_{A_n,\sigma_n}(x_k)$.  This requires a numerical optimization procedure, such as gradient ascent.  The problem is concave when $\log\psi$ is concave.
   (ii)  Minimize (20.21) with respect to $A, \sigma^2$, yielding

$$A_{n+1} = \left(\sum_{k=1}^{N}x_k b_k^T\right)\left(\sum_{k=1}^{N}S_k\right)^{-1}$$

with $b_k = \hat{y}_{A_n,\sigma_n}(x_k)$, $S_k = \hat{y}_{A_n,\sigma_n}(x_k)\hat{y}_{A_n,\sigma_n}(x_k)^T$, and

$$\sigma_{n+1}^2 = \frac{1}{Nd}\sum_{k=1}^{N}\left|x_k - A\hat{y}_{A_n,\sigma_n}\right|^2.$$

   (3)  Stop if the variation of the parameter is below a tolerance level.  Otherwise, iterate to the next step.

---

Once $A$ and $\sigma^2$ have been estimated, the $y$ components associated to a new observation $x$ can be estimated by $\hat{y}_{A,\sigma}(x)$, therefore minimizing

$$\frac{1}{2\sigma^2}\left|x_k - Ay\right|^2 + \sum_{j=1}^{p}\log\psi(y^{(j)}),$$

yielding the map estimate, the same convex optimization problem as in step (1) above. Now we can see how the method takes from both PCA and ICA: the columns

of $A$, $a_1, \ldots, a_p$ can be considered as $p$ principal directions, and are fixed after learning; they are not orthonormal, and do not satisfy the nesting properties of PCA (that those $p$ contain those for $p-1$). The coordinates of $x$ with respect to this basis is not a projection, as would be provided by PCA, but the result of a penalized estimation problem. The penalty associated to the logistic case is

$$\log \psi(y^{(j)}) = \log 2 - 2\log(e^{y^{(j)}} + e^{-y^{(j)}}).$$

This distribution with "exponential tails" has the interest of allowing large values of $y^{(j)}$, which generally entails *sparse decompositions*, in which $y$ has a few large coefficients, and many zeros.

As an alternative to the mode approximation of the EM, which may lead to biased estimators, one may use the SAEM algorithm (section 16.4.3), as proposed in Allassonniere and Younes [3]. Recall that the EM algorithm replaces the parameters $A_0, \sigma_0^2$ by minimizers of

$$\frac{Nd}{2}\log(\sigma^2) + \frac{1}{2\sigma^2}\sum_{k=1}^{N} E_{A_0,\sigma_0}(|x_k - AY|^2 | X = x_k)$$

$$= \frac{Nd}{2}\log(\sigma^2) + \frac{1}{2\sigma^2}\sum_{k=1}^{N}|x_k|^2 - \frac{1}{\sigma^2}\sum_{k=1}^{N}x_k^T A b_k + \frac{1}{2\sigma^2}\sum_{k=1}^{N}\text{trace}(A^T A S_k),$$

where the computation of $b_k^{(j)} = E_{A_0,\sigma_0}(Y^{(j)}|X = x_k)$ and $s_k(i,j) = E_{A_0,\sigma_0}(Y^{(i)}Y^{(j)}|X = x_k)$ was the challenging issue. In the SAEM algorithm, the statistics $b_k$ and $S_k$ are part of a stochastic approximation scheme, and are estimated in parallel with EM updates as follows.

---

**Algorithm 20.3 (SAEM for probabilistic ICA)**
Initialize the algorithm with parameters $A$, $\sigma^2$. Define a sequence of decreasing steps, $\gamma_t$.

Let, for $k = 1, \ldots, N$, $b_k = 0$ and $S_k = \text{Id}$. Iterate the following steps.

(1) For $k = 1, \ldots, N$, sample $y_k$ according to the conditional distribution of $Y$ given $X = x_k$, using the current parameters $A$ and $\sigma^2$.

(2) Update $b_k$ and $S_k$, letting (assuming step $t$ of the algorithm)

$$\begin{cases} b_k \to b_k + \gamma_t(Y_k - b_k) \\ S_k \to S_k + \gamma_t(Y_k Y_k^T - S_k) \end{cases}$$

(3) Replace $A$ and $\sigma^2$ by

$$A = \left(\sum_{k=1}^{N} x_k b_k^T\right)\left(\sum_{k=1}^{N} S_k\right)^{-1}$$

and

$$\sigma^2 = \frac{1}{Nd}\sum_{k=1}^{N}\left|x_k - A\hat{y}_{A_0,\sigma_0}\right|^2.$$

---

The parameter $\gamma_t$ should be decreasing with $t$, typically so that $\sum_t \gamma_t = +\infty$ and $\sum_t \gamma_t^2 < \infty$ (e.g., $\gamma_t \propto 1/t$). One way to sample from $Y_k$ is to uses a rejection scheme, iterating the procedure which samples $y$ according to the prior and accepts the result with probability $M\exp(-|x_k - Ay|^2/2\sigma^2)$ until acceptance. Here $M$ must be chosen so that $M\max_y \exp(-|x_k - Ay|^2/2\sigma^2) \le 1$ (e.g., $M = 1$).

This method will work for small $p$, but for large $p$, the probability of acceptance may be very small. In such cases, $Y_k$ can be sampled changing one component at a time using a Metropolis-Hastings scheme. If component $j$ is updated, this scheme samples a new value of $y$ (call it $y'$) by changing only $y^{(j)}$ according to the prior distribution $\psi$ and accept the change with probability

$$\min\left(1, \frac{\exp(-|x_k - Ay'|^2/2\sigma^2)}{\exp(-|x_k - Ay|^2/2\sigma^2)}\right).$$

## 20.7   Non-negative matrix factorization

In this section, we consider factor analysis methods that approximate a random variable $X$ in the form $X = \sum_{j=1}^{p} a^{(j)} Y^{(j)}$ with the constraint that the scalars $a^{(1)}, \ldots, a^{(p)} \in \mathbb{R}$ and the vectors $Y^{(1)}, \ldots, Y^{(p)} \in \mathbb{R}^d$ are respectively non-negative and with non-negative entries. This model makes sense, for example, when $X$ represents the total multivariate production (e.g., in terms of number of molecules of various types) resulting of several chemical reactions that operate together. Another application is when $X$ is a list of preference scores associated with a person for, say, books or movies, and each person is modeled as a positive linear combination of $p$ "typical scorers," represented by the vector $Y^{(j)}$ for $j = 1, \ldots, p$.

When training data $(x_1, \ldots, x_N)$ is observed and stacked in an $N$ by $d$ matrix $\mathcal{X}$, the decomposition can be summarized for all observations together in the matrix form

$$\mathcal{X} = \mathcal{A}Y^T$$

where $\mathcal{A}$ is $N$ by $p$ and provides the coefficients $a_k^{(j)}$ associated with each observation and $Y = [y^{(1)}, \ldots, y^{(p)}]$ is $d$ by $p$ and provides the $p$ typical profiles. The matrices $\mathcal{A}$ and $Y$ are unknown and their estimation subject to the constraint of having non-negative components represent the non-negative matrix factorization (NMF) problem.

NMF is often implemented by solving the constrained optimization problem of minimizing $|\mathcal{X} - \mathcal{A} Y^T|^2$ subject to $\mathcal{A}$ and $Y$ having non-negative entries. This problem is non-convex in general but the sub-problems of optimizing either $\mathcal{A}$ or $Y$ when the other matrix is fixed are simple quadratic programs.

This suggests using an alternating minimization method, iterating steps in which $\mathcal{A}$ is updated with $Y$ fixed, followed by an update of $Y$ with $\mathcal{A}$ fixed. However, solving a full quadratic program at each step would be computationally prohibitive with large datasets, and simpler update rules have been suggested, updating each matrix in turn with a guarantee of reducing the objective function.

If $Y$ is considered as fixed and $\mathcal{A}$ is the free variable, we have

$$|\mathcal{X} - \mathcal{A} Y^T|^2 = |\mathcal{X}|^2 - 2\mathrm{trace}(\mathcal{X}^T \mathcal{A} Y^T) + \mathrm{trace}(\mathcal{A} Y^T Y \mathcal{A}^T)$$
$$= \mathrm{trace}(\mathcal{A}^T \mathcal{A} Y^T Y) - 2\mathrm{trace}(\mathcal{A}^T (\mathcal{X} Y)) + |\mathcal{X}|^2.$$

The next lemma will provide update steps for $\mathcal{A}$.

**Lemma 20.14** *Let $M$ be an $n$ by $n$ symmetric matrix and $b \in \mathbb{R}^n$, both assumed to have non-negative entries. Let $u \in \mathbb{R}^n$, also with non-negative coefficients, and let*

$$v^{(i)} = u^{(i)} \left( \frac{b^{(i)}}{\sum_{j=1}^d m(i,j) u^{(j)}} \right).$$

*Then*

$$v^T M v - 2b^T v \le u^T M u - 2b^T u.$$

*Moreover, $v = u$ if and only if $u$ minimizes $u^T M u - 2b^T u$ subject to $u^{(i)} = 0$, $i = 1, \ldots, n$.*

PROOF Let $F(u) = u^T M u - 2b^T u$. We look for $v^{(i)} = \beta^{(i)} u^{(i)}$ with $\beta^{(i)} \ge 0$ such that $F(v) \le F(u)$. We have

$$F(v) = \sum_{i,j=1}^n \beta^{(i)} \beta^{(j)} u^{(i)} u^{(j)} m(i,j) - 2 \sum_{i=1}^n b^{(i)} \beta^{(i)}$$

$$\le \frac{1}{2} \sum_{i,j=1}^n ((\beta^{(i)})^2 + (\beta^{(j)})^2) u^{(i)} u^{(j)} m(i,j) - 2 \sum_{i=1}^n b^{(i)} \beta^{(i)} u^{(i)}$$

$$= \sum_{i,j=1}^n (\beta^{(i)})^2 u^{(i)} u^{(j)} m(i,j) - 2 \sum_{i=1}^n b^{(i)} \beta^{(i)} u^{(i)}$$

When $\beta = \mathbb{1}_n$, this upper-bound is equal to $F(u)$. So, if we choose $\beta$ minimizing the upper-bound, we will indeed find $v$ such that $F(v) \leq F(u)$. Rewriting the upper-bound as

$$\sum_{i=1}^{n} u^{(i)} \left( (\beta^{(i)})^2 \left( \sum_{j=1}^{n} m(i,j) u^{(j)} \right) - 2 b^{(i)} \beta^{(i)} \right)$$

we see that $\beta^{(i)} = b^{(i)} / \sum_{j=1}^{n} m(i,j) u^{(j)}$ provides such a minimizer, which proves the first statement of the lemma. For the second statement, we have $v^{(i)} = u^{(i)}$ if and only if $u^{(i)} = 0$ or $\sum_{j=1}^{n} m(i,j) u^{(j)} = b^{(i)}$, and one directly checks that these are exactly the KKT conditions for a minimizer of $F$ over vectors with non-negative entries. ∎

To apply the lemma to the minimization in $\mathcal{A}$, let $M : \mathcal{A} \mapsto \mathcal{A} Y^T Y$ and $b = \mathcal{X} Y$ (we are working in the linear space of $N$ by $p$ matrices). Then the update

$$a_k^{(i)} \mapsto a_k^{(i)} \frac{(\mathcal{X} Y)(i,k)}{(\mathcal{A} Y^T Y)(i.k)}$$

decreases the objective function.

Similarly, applying the lemma with the operator $Y \mapsto Y \mathcal{A}^T \mathcal{A}$ and $b = \mathcal{X}^T \mathcal{A}$ gives the update for $Y$, namely

$$y_j^{(i)} \mapsto y_j^{(i)} \frac{(\mathcal{X}^T \mathcal{A})(i,j)}{(Y \mathcal{A}^T \mathcal{A})(i,j)}.$$

We have therefore obtained the following algorithm.

---

**Algorithm 20.4 (NMF, quadratic cost)**

1. Fix $p > 0$ and let $\mathcal{X}$ be the $N$ by $d$ matrix containing the observed data. Initialize the procedure with matrices $A$ and $Y$, respectively of size $N$ by $p$ and $d$ by $p$, with positive coefficients.

2. At a given stage of the algorithm, let $\mathcal{A}$ and $Y$ be the current matrices providing an approximate decomposition of $\mathcal{X}$.

3. For the next step, let $\tilde{A}$ be the matrix with coefficients

$$\tilde{a}_k^{(i)} = a_k^{(i)} \frac{(\mathcal{X} Y)(i,k)}{(\mathcal{A} Y^T Y)(i,k)}$$

and $\tilde{Y}$ the matrix with coefficients

$$\tilde{y}_j^{(i)} = y_j^{(i)} \frac{(\mathcal{X}^T \tilde{A})(i,j)}{(Y \tilde{A}^T \tilde{A})(i,j)}.$$

4. Replace $\mathcal{A}$ by $\tilde{\mathcal{A}}$ and $Y$ by $\tilde{Y}$, iterating until numerical convergence.

---

An alternative version of the method has been proposed, where the objective function is $\Phi(\mathcal{A}Y^T)$, where, for an $N$ by $d$ matrix $\mathcal{Z} = [z_1, \ldots, z_N]^T$,

$$\Phi(\mathcal{Z}) = \sum_{k=1}^{N} \sum_{i=1}^{d} (z_k^{(i)} - x_k^{(i)} \log z_k^{(i)})$$

which is indeed minimal for $\mathcal{Z} = \mathcal{X}$. We state and prove a second lemma that will allow us to address this problem.

**Lemma 20.15** *Let $M$ be an $n$ by $q$ matrix and $x \in \mathbb{R}^n$, $b \in \mathbb{R}^q$, all assumed to have positive entries. For $u \in (0, +\infty)^q$, define*

$$F(u) = \sum_{j=1}^{q} b^{(j)} u^{(j)} - \sum_{i=1}^{n} x^{(i)} \log \sum_{j=1}^{q} m(i,j) u^{(j)}.$$

*Define $v \in (0, +\infty)^q$ by*

$$v^{(j)} = u^{(j)} \left( \frac{\sum_{i=1}^{n} m(i,j) x^{(i)}/\alpha^{(i)}}{b^{(j)}} \right)$$

*with $\alpha^{(i)} = \sum_{k=1}^{q} m(i,k) u^{(k)}$. Then $F(v) \le F(u)$. Moreover, $v = u$ if and only if $u$ minimizes $F$ subject to $u^{(i)} \ge 0, i = 1, \ldots, n$.*

PROOF Introduce a variable $\beta^{(j)} > 0$ for $j = 1, \ldots, q$ an let $w^{(j)} = u^{(j)} \beta^{(j)}$. Then

$$F(w) = \sum_{j=1}^{q} b^{(j)} u^{(j)} \beta^{(j)} - \sum_{i=1}^{n} x^{(i)} \log \sum_{j=1}^{q} m(i,j) u^{(j)} \beta^{(j)}$$

$$= \sum_{j=1}^{q} b^{(j)} u^{(j)} \beta^{(j)} - \sum_{i=1}^{n} x^{(i)} \log \frac{\sum_{j=1}^{q} m(i,j) u^{(j)} \beta^{(j)}}{\sum_{j=1}^{q} m(i,j) u^{(j)}} - \sum_{i=1}^{n} x^{(i)} \log \sum_{j=1}^{q} m(i,j) u^{(j)}$$

Let $\rho(i,j) = m(i,j) u^{(j)}/\alpha^{(i)}$. Since the logarithm is concave, we have

$$\log \sum_{j=1}^{q} \rho(i,j) \beta^{(j)} \ge \sum_{j=1}^{q} \rho(i,j) \log \beta^{(j)}$$

so that

$$F(w) \le \sum_{j=1}^{q} b^{(j)} u^{(l)} j \beta^{(j)} - \sum_{i=1}^{n} \sum_{j=1}^{q} x^{(i)} \rho(i,j) \log \beta^{(j)} - \sum_{i=1}^{n} x^{(i)} \log \sum_{j=1}^{q} m(i,j) u^{(j)}.$$

The upper bound with $\beta^{(j)} \equiv 1$ gives $F(u)$, so minimizing this expression in $\beta$ will give $F(w) \leq F(u)$. This minimization is straightforward and gives

$$\beta^{(j)} = \frac{\sum_{i=1}^{n} x^{(i)} \rho(i,j)}{b^{(j)} u^{(j)}} = \frac{\sum_{i=1}^{n} m(i,j) x^{(i)} / \alpha^{(i)}}{b^{(j)}}$$

and the optimal $w$ is the vector $v$ provided in the lemma. Finally, one checks that $v = u$ if and only if $u$ satisfies the KKT conditions for the considered problem.  ∎

We can now apply this lemma to derive update rules for $\mathcal{Y}$ and $\mathcal{A}$, where the objective is

$$\sum_{k=1}^{N} \sum_{i=1}^{d} \sum_{j=1}^{p} y_j^{(i)} a_k^{(j)} - \sum_{k=1}^{N} \sum_{i=1}^{d} x_k^{(i)} \log \sum_{j=1}^{p} y_j^{(i)} a_k^{(j)}.$$

Starting with the minimization in $\mathcal{A}$, we apply the lemma to each index $k$ separately, taking $n = d$ and $q = p$, with $b^{(j)} = \sum_{i=1}^{d} y_j^{(i)}$ and $m(i,j) = y_i^{(j)}$. Then the update is

$$a_k(j) \mapsto a_k(j) \frac{\sum_{i=1}^{d} x_k^{(i)} y_j^{(i)} / \alpha_k^{(i)}}{\sum_{i=1}^{d} y_j^{(i)}}$$

with $\alpha_k^{(i)} = \sum_{j=1}^{p} y_j^{(i)} a_k^{(j)}$.

For $Y$, we can work with fixed $i$ and apply the lemma with $n = N$, $q = p$, $b^{(j)} = \sum_{k=1}^{N} a_k^{(j)}$ and $m(k,j) = a_k^{(j)}$. This gives the update:

$$y_j^{(i)} \mapsto y_j^{(i)} \frac{\sum_{k=1}^{N} x_k^{(i)} a_k^{(j)} / \alpha_k^{(i)}}{\sum_{k=1}^{N} a_k^{(j)}},$$

still with $\alpha_k^{(i)} = \sum_{j=1}^{p} y_j^{(i)} a_k^{(j)}$.

We summarize this in our second algorithm for NMF.

---

**Algorithm 20.5 (NMF, logarithmic cost)**

1. Fix $p > 0$ and let $\mathcal{X}$ be the $N$ by $d$ matrix containing the observed data.

2. Initialize the procedure with matrices $Y$ and $\mathcal{A}$, respectively of size $N$ by $p$ and $d$ by $p$, with positive coefficients.

3. At a given stage of the algorithm, let $\mathcal{A}$ and $Y$ be the current matrices decomposing $\mathcal{X}$.

4. Let $\tilde{A}$ be the matrix with coefficients

$$\tilde{a}_k^{(j)} = a_k^{(j)} \frac{\sum_{i=1}^d x_k^{(i)} y_j^{(i)} / \alpha_k^{(i)}}{\sum_{i=1}^d y_j^{(i)}}$$

with $\alpha_k^{(i)} = \sum_{j=1}^p y_j^{(i)} a_k^{(j)}$.

5. Let $\tilde{Y}$ the matrix with coefficients

$$\tilde{y}_j^{(i)} = y_j^{(i)} \frac{\sum_{k=1}^N x_k^{(i)} \tilde{a}_k^{(j)} / \tilde{\alpha}_k^{(i)}}{\sum_{j=1}^p \tilde{a}_k^{(j)}}$$

with $\tilde{\alpha}_k^{(i)} = \sum_{j=1}^p y_j^{(i)} \tilde{a}_k^{(j)}$.

6. Replace $\mathcal{A}$ by $\tilde{A}$ and $Y$ by $\tilde{Y}$, iterating until numerical convergence.

---

## 20.8 Variational Autoencoders

Variational autoencoders, which were described in section 18.2.2, can be ineterpreted as a non-linear factor model in which $X = g(\theta, Y) + \epsilon$ where $\epsilon$ is a centered Gaussian noise with covariance matrix $Q$ and $Y \in \mathbb{R}^p$ has a known probability distribution, such as $Y \sim \mathcal{N}(0, \mathrm{Id}_{\mathbb{R}^p})$. In this framework, the conditional distribution of $Y$ given $X = x$ was approximated as a Gaussian distribution with mean $\mu(x, w)$ and covariance matrix $S(x, w)^2$. The implementation in Kingma and Welling [103, 104] use neural networks for the three functions $g$, $\mu$ and $S$.

## 20.9 Bayesian factor analysis and Poisson point processes

### 20.9.1 A feature selection model

The expectation in many factor models is that individual observations are obtained by mixing pure categories, or topics, and represented as a weighted sum or linear combination of a small number of uncorrelated or independent variables. Denote $p$ the number of possible categories, which, in this section, can be assumed to be quite large.

We will assume that each observation randomly selects a small number among these categories before combining them. Let us consider (as an example) the following model.

- The observations $X_1, \ldots, X_N$ take the form of a probabilistic ICA model

$$X_k = \sum_{j=1}^{p} a_k(j) b_k(j) Y^{(j)} + \sigma R_k,$$

where:

- $R_k$ follows a standard Gaussian distribution,

- $a_k(1), \ldots, a_k(p)$ are independent with $a_k(j) \sim \mathcal{N}(m_j, \tau_j^2)$,

- $b_k(1), \ldots, b_k(p)$ are independent and follow a Bernoulli distribution with parameter $\pi_j$,

- $Y^{(1)}, \ldots, Y^{(p)}$ are independent standard Gaussian random variables.

- $\sigma^2$ follows an inverse gamma distribution with parameters $\alpha_0, \beta_0$.

- $\tau_1^2, \ldots, \tau_p^2$ follow independent inverse gamma distributions with parameters $\alpha_1, \beta_1$.

- $m_j$ follow a Gaussian $\mathcal{N}(0, \rho^2)$ and,

- $\pi_j$ follow a beta distribution with parameters $(u, v)$.

The priors are, as usual, chosen so that the computation of posterior distributions is easy, i.e., they are conjugate priors. The observed data is therefore obtained by selecting components $Y_j$ with probability $\pi_j$ and weighted with a Gaussian random coefficient, then added before introducing noise.

Let $n_j = \sum_{k=1}^{N} b_k(j)$. Ignoring constant factors, the joint likelihood of all variables together is proportional to:

$$L \propto \sigma^{-Nd} \exp\left(-\frac{1}{2\sigma^2} \sum_{k=1}^{N} |X_k - \sum_{j=1}^{p} a_k(j) b_k(j) Y^{(j)}|^2\right)$$

$$\prod_{j=1}^{p} \left(\tau_j^{-N} \exp\left(-\frac{1}{2\tau_j^2} \sum_{k=1}^{N} (a_k(j) - m_j)^2\right)\right) \exp\left(-\frac{1}{2\rho^2} \sum_{j=1}^{p} m_j^2\right)$$

$$\prod_{j=1}^{p} \left(\pi_j^{n_j} (1 - \pi_j)^{N - n_j}\right) \prod_{j=1}^{p} \left((\tau_j^2)^{-\alpha_1 - 1} \exp(-\beta_1 / \tau_j^2)\right)$$

$$(\sigma^2)^{\alpha_0 - 1} \exp(-\beta_0 / \sigma^2) \prod_{j=1}^{p} \left(\pi_j^{u-1} (1 - \pi_j)^{v-1}\right) \exp\left(-\frac{1}{2} \sum_{i=1}^{p} |Y^{(i)}|^2\right)$$

In spite of the complexity of this expression, it is relatively straightforward (by considering each variable in isolation) to see that

- The conditional distribution of $\sigma^2, \tau_1^2, \ldots, \tau_p^2$ given all other variables remains a product of inverse gamma distributions.

- The conditional distribution of $Y^{(1)}, \ldots, Y^{(p)}$ given the other variables is Gaussian.

- The conditional distribution of $\pi_1, \ldots, \pi_p$ given the other variables is a product of beta distributions.

- The conditional distribution of $m_1, \ldots, m_p$ given the other variables remain independent Gaussian.

- The posterior distribution of $a_1, \ldots, a_N$ (considered as $p$-dimensional vectors) given the other variables is a product of independent Gaussian (but the components $a_k(j)$, $j = 1, \ldots, p$ are correlated).

- For the posterior distribution given the other variables, $b_1, \ldots, b_N$ (considered as $p$-dimensional vectors) are independent. The components of each $b_k$ are not independent but each $b_k(j)$ being a binary variable follows a Bernoulli distribution given the other ones.

These remarks provide the basis of a Gibbs sampling algorithm for the simulation of the posterior distribution of all unobserved variables (the computation of the parameters of each of the conditional distribution above requires some work, of course, and these details are left to the reader). This simulation does not explicitly provide a matrix factorization of the data (in the sense of a single matrix $\mathcal{A}$ such that $\mathcal{X} = \mathcal{A}Y$, as considered in the previous section), but a probability distribution on such matrices, expressed as $\mathcal{A}(k, j) = a_k(j)b_k(j)$. One can however use the average of the matrices obtained through the simulation for this purpose. Additional information can be obtained through this simulation. For example, the expectation of $b_k(j)$ provides a measure of proximity of observation $k$ to category $j$.

### 20.9.2  Non-negative and count variables

**Poisson factor analysis.**  Many variations can be made on the previous construction. When the observations are non-negative, for example, an additive Gaussian noise may not be well adapted. Alternative models should model the conditional distribution of $X$ given $a$, $b$ and $Y$ as a distribution over non-negative numbers with mean $(a \odot b)^T Y$ (for example a gamma distribution with appropriate parameters). The posterior sampling generally is more challenging in this case because simple conjugate priors are not always available.

An important special case is when $X$ is a count variable taking values in the set of non-negative integers. In this case (starting with a model without feature selection), modeling $X$ as a Poisson variable with mean $a(1)Y^{(1)} + \cdots + a(p)Y^{(p)}$ leads to tractable computations, once it is noticed that $X$ can be seen as a sum of random variables

$Z^{[1]}, \ldots, Z^{[p]}$ where $Z^{[i]}$ follows a Poisson distribution with parameter $a(i)Y^{(i)}$. This suggests introducing new latent variables $(Z^{[1]}, \ldots, Z^{[p]})$, which are not observed but follow, conditionally to their sum, which is $X$ and is observed, a multinomial distribution with parameters $X, q_1, \ldots, q_p$, with $q_i = a(i)Y^{(i)}/(\sum_{j=1}^{p} a(j)Y^{(j)})$.

This provides what is referred to as a Poisson factor analysis (PFA). As an example, consider a Bayesian approach where, for the prior distribution, $a(1), \ldots, a(p)$ are independent and follow as a gamma distribution with parameters $\alpha_0$ and $\beta_0$, and $Y^{(1)}, \ldots, Y^{(p)}$ are independent, exponentially distributed with parameter 1. The joint likelihood of all data then is (up to constant factors):

$$L \propto \exp\left(-\sum_{k=1}^{N}(a_k(1)Y^{(1)} + \cdots + a_k(p)Y^{(p)})\right)\left(\prod_{k=1}^{N}\prod_{i=1}^{p}\frac{(a_k(i)Y^{(i)})^{z_k^{[i]}}}{z_k^{[i]}!}\right)$$

$$\left(\prod_{k=1}^{N}\prod_{i=1}^{p} a_k(i)^{\alpha-1}\right)\exp\left(-\beta\sum_{k=1}^{N}\sum_{i=1}^{p} a_k(i)\right)\exp\left(-\sum_{i=1}^{p} Y^{(i)}\right).$$

This is the GaP (For Gamma-Poisson) model introduced in Canny [50]. The conditional distribution of the variables $(a_k(i))$ given $(Z_k^{[i]})$ and $(Y_i)$ are independent and gamma-distributed, and so are $(Y^{(i)})$ given the other variables. Finally, for each $k$, the family $(Z_k^{[1]}, \ldots, Z_k^{[p]})$ follows a multinomial distribution conditionally to their sum, $X_k$, and the rest of the variables, and these variables are conditionally independent across $k$.

**GaP with feature selection**    One can include a feature selection step in this model by introducing binary variables $b(1), \ldots, b(p)$, with selection probabilities $\pi_1, \ldots, \pi_p$, with a Beta$(u, v)$ prior distribution on $\pi_i$. Doing so, the likelihood of the extended model is:

$$L \propto \exp\left(-\sum_{k=1}^{N}(a_k(1)b_k(1)Y^{(1)} + \cdots + a_k(p)b_k(p)Y^{(p)})\right)\left(\prod_{k=1}^{N}\prod_{i=1}^{p}\frac{(a_k(i)b_k(i)Y^{(i)})^{z_k^{[i]}}}{z_k^{[i]}!}\right)$$

$$\left(\prod_{k=1}^{N}\prod_{i=1}^{p} a_k(i)^{\alpha-1}\right)\exp\left(-\beta\sum_{k=1}^{N}\sum_{i=1}^{p} a_k(i)\right)\exp\left(-\sum_{i=1}^{p} Y^{(i)}\right)$$

$$\prod_{j=1}^{p}\left(\pi_j^{n_j}(1-\pi_j)^{N-n_j}\right)\prod_{j=1}^{p}\left(\pi_j^{u-1}(1-\pi_j)^{v-1}\right).$$

where, as before, $n_j = \sum_{k=1}^{N} b_k(j)$. The conditional distribution of $\pi_1, \ldots, \pi_p$ given the other variables is therefore still that of a family of independent beta-distributed variables. The binary variables $b_k(1), \ldots, b_k(p)$ are also conditionally independent given the other variables, with $b_k(i) = 1$ with probability one if $z_k^{[i]} > 0$ and with probability $\pi_j \exp(-a_k(j)Y^{(j)})$ if $z_k^{[j]} = 0$.

### 20.9.3   Feature assignment model

The previous models assumed that $p$ features were available, modeled as $p$ random variables with some prior distribution, and that each observation picks a subset of them, drawing feature $j$ with probability $\pi_j$. We denoted by $b_k(j)$ the binary variable indicating whether feature $j$ was selected for observation $k$, and $n_j$ was the number of times that feature was selected. Finally, we modeled $\pi_j$ as a beta variable with parameters $u$ and $v$.

One can compute, using this model, the probability distribution of of the feature selection variables, $\boldsymbol{b} = (b_k(j), j = 1,\dots,p, k = 1,\dots,N)$. From the model definition, the probability of observing such a configuration is given by

$$
Q(\boldsymbol{b}) = \frac{\Gamma(u+v)^p}{\Gamma(u)^p \Gamma(v)^p} \int \prod_{j=1}^p \pi_j^{n_j+u-1} (1-\pi_j)^{N-n_j+v-1} d\pi_1 \dots d\pi_p
$$

$$
= \prod_{j=1}^p \frac{\Gamma(u+v)\Gamma(u+n_j)\Gamma(v+N-n_j)}{\Gamma(u)\Gamma(v)\Gamma(u+v+N)}
$$

$$
= \prod_{j=1}^p \frac{u(u+1)\cdots(u+n_j-1)v(v+1)\cdots(v+N-n_j-1)}{(u+v)(u+v+1)\cdots(u+v+N-1)}
$$

Denote by $n_{jk} = \sum_{l=1}^{k-1} b_l(j)$ the number of observations with index less than $k$ that pick feature $j$. Using this notation, we can write, using the fact that

$$
u(u+1)\cdots(u+n_j-1) = \prod_{k=1}^N (u+n_{jk})^{b_k(j)}
$$

and a similar identity for $v(v+1)\cdots(v+N-n_j-1)$,

$$
Q(\boldsymbol{b}) = \prod_{j=1}^p \prod_{k=1}^N \frac{(u+n_{jk})^{b_k(j)}(v+k-1-n_{jk})^{1-b_k(j)}}{u+v+k-1}
$$

$$
= \prod_{k=1}^N \prod_{j=1}^p \left(\frac{u+n_{jk}}{u+v+k-1}\right)^{b_k(j)} \left(\frac{v+k-1-n_{jk}}{u+v+k-1}\right)^{1-b_k(j)}.
$$

Using this last equation, we can interpret the probability $Q$ as resulting from a progressive feature assignment process. The first observation, $k = 1$, for which $n_{jk} = 0$ for all $j$, chooses each feature with probability $u/(u+v)$. When reaching observation $k$, feature $j$ is chosen with probability $(u+n_{jk})/(u+v+k-1)$. At all steps, features are chosen independently from each other.

Let $F_k$ be the set of features assigned to observation $k$, i.e., $F_k = \{j : b_k(j) = 1\}$ and

$$G_k = F_k \setminus \bigcup_{l=1}^{k-1} F_l$$

be the set of features used in observation $k$ but in no previous observation. Let $C_k = F_k \setminus G_k$ and $U_k = G_1 \cap \cdots \cap G_{k-1}$. Instead of considering configurations $b = (b_k(j), i = 1, \ldots, d, k = 1, \ldots, N)$ we may alternatively consider the family of sets $S = (G_k, C_k, 1 \le k \le N)$. Such a family must satisfy the property that the sets $G_k$ and $C_k$ are non-intersecting, $C_k \subset U_k$ and $G_l \cap G_k = \emptyset$ for $l < k$. It provide a unique configuration $b$ by letting $b_k(j) = 1$ if and only if $j \in G_k \cup C_k$. We will let, in the following, $q_k = |G_k|$ and $p_k = |U_k|$. The probability $Q(b)$ can be re-expressed in terms of $S$, letting (with some abuse of notation)

$$Q(S) = \prod_{k=1}^{N} \left( \left( \frac{u}{u+v+k-1} \right)^{q_k} \left( \frac{v+k-1}{u+v+k-1} \right)^{p-p_{k+1}} \right.$$
$$\left. \prod_{j \in U_k} \left( \frac{u+n_{jk}}{u+v+k-1} \right)^{\mathbf{1}_{j \in C_k}} \left( \frac{v+k-1-n_{jk}}{u+v+k-1} \right)^{\mathbf{1}_{j \notin C_k}} \right).$$

Let $S_k = (G_l, C_l, l \le k)$. Then the expression of $Q$ shows that, conditionally to $S_{k-1}$, $G_k$ and $C_k$ are independent. Elements in $C_k$ are chosen independently for each feature $j \in U_k$ with probability $(u+n_{jk})/(u+v+k-1)$. Moreover, the conditional distribution of $q_k$ given $S_{k-1}$ is proportional to

$$\left( \frac{u}{u+v+k-1} \right)^{q_k} \left( \frac{v+k-1}{u+v+k-1} \right)^{p-p_k-q_k}$$

i.e., it is a binomial distribution with parameters $p - p_k$ and $u/(u+v+k-1)$. Finally, given $s_{k-1}$ and $q_k$, the distribution of $G_k$ is uniform among all $\binom{p-p_k}{q_k}$ subsets of

$$\{1, \ldots, p\} \setminus (G_1 \cup \cdots \cup G_{k-1})$$

with cardinality $q_k$.

If there is no special meaning in the feature label, which is the case in our discussion of prior models in which all features are sampled independently with the same distribution, we may identify configurations that can be deduced from each other by relabeling (note that relabeling features does not change the value of $Q$).

Call a configuration normal if $G_k = \{p_k + 1, \ldots, p_{k+1}\}$. Given $S$, it is always possible to relabel the features with a permutation $\sigma$ so that, for each $k$, $\sigma(G_k) = \{p_k + 1, \ldots, p_{k+1}\}$. There are, in fact, $q_1! \ldots q_N!$ such permutations. We can complete the process generating $S$ by adding at the end a transformation into a normal configuration

(picking uniformly at random one of the possible ones). The probability of a normal configuration $S$ obtained through this process is (using a simple counting argument)

$$Q(S) = \prod_{k=1}^{N} \left( \binom{p-p_k}{q_k} \left( \frac{u}{u+v+k-1} \right)^{q_k} \left( \frac{v+k-1}{u+v+k-1} \right)^{p-p_{k+1}} \right.$$
$$\left. \prod_{j \in U_k} \left( \frac{u+n_{jk}}{u+v+k-1} \right)^{\mathbf{1}_{j \in C_k}} \left( \frac{v+k-1-n_{jk}}{u+v+k-1} \right)^{\mathbf{1}_{j \notin C_k}} \right),$$

This provides a new incremental procedure that directly samples normalized assignments. First let $q_1$ follow a binomial distribution $\mathrm{bin}(p, u/(u+v))$ and assign the first observation to features 1 to $q_1$. Assume that $p_k$ labels have been created before step $k$. Then select for observation $k$ some of the already labeled features, label $j$ being selected with probability $(u+n_{jk})/(u+v+k-1)$ as above. Finally, add $q_k$ new features where $q_k$ follows a binomial distribution $\mathrm{bin}(p-p_k, u/(u+v+k-1))$.

This discussion is clearly reminiscent of the one that was made in section 19.7.3 leading to the Polya urn process, and we want here also to let $p$ tend to infinity (with fixed $N$) with proper choices of $u$ and $v$ as functions of $p$ in the expression above. Choose two positive numbers $c$ and $\gamma$ and let $u = c\gamma/p$ and $v = c - u$. Note that, with the incremental simulation process that we just described, the conditional expectation of the next number of labels, $p_{k+1}$ given the current one, $p_k$ is

$$E(p_{k+1}|p_k) = \frac{(p-p_k)u}{u+v+k-1} + p_k = \frac{c\gamma}{c+k-1} + \left( 1 - \frac{c\gamma}{p(c+k-1)} \right) p_k \le \frac{c\gamma}{c+k-1} + p_k$$

Taking expectations on both sides, we get

$$E(p_{k+1}) \le \sum_{l=1}^{k} \frac{c\gamma}{c+l-1} \le \sum_{l=1}^{N} \frac{c\gamma}{c+l-1}$$

so that this expectation is bounded independently of $k$. This shows in particular that $p_k/p$ tends to 0 in probability (just applying Markov's inequality) and that the binomial distribution $\mathrm{bin}(p-p_k, u/(u+v+k-1))$ can be approximated by a Poisson distribution with parameter $c\gamma/(c+k-1)$.

So, when $p \to \infty$, we obtain the following incremental simulation process for the feature labels, that we combine with the actual simulation of the features, assumed to follow a prior distribution with p.d.f. $\psi$. This process is called the Indian buffet process in the literature, the analogy being that a buffet offers an infinite variety of dishes, and each observation is a customer who tastes a finite number of them.

**Algorithm 20.6 (Indian buffet process)**
  1.  Initialization:

       (i) Sample an integer $q_1$ according to a Poisson distribution with parameter $\gamma$.

       (ii) Sample features $y^{(1)},\ldots,y^{(q_1)}$ according to $\psi$.

       (iii) Assign these features to observation 1, and let $n_{2,j} = 1$ for $j = 1,\ldots,q_1$.

   2. Assume that observations 1 to $k-1$ have been obtained, with $p_k$ features $y^{(1)},\ldots,y^{(p_k)}$ such that the $j$th feature has been chosen $n_{k,j}$ times.

       (i) For $j = 1,\ldots,p_k$, assign feature $j$ to sample $k$ with probability $\frac{n_{k,j}}{c+k-1}$. If $j$ is selected, let $n_{k+1,j} = n_{k,j} + 1$, otherwise let $n_{k+1,j} = n_{k,j}$.

       (ii) Sample an integer $q_k$ according to a Poisson distribution with parameter $\frac{c\gamma}{c+k-1}$ and let $p_{k+1} = p_k + q_k$.

       (iii) Sample features $y^{(p_k+1)},\ldots,y^{(p_{k+1})}$ according to $\psi$.

       (iv) Assign these features to observation k, and let $n_{k+1,j} = 1$ for $j = p_k+1,\ldots,p_k$.

   3. If $k = N$, stop, otherwise replace $k$ by $k+1$ and return to Step 2.

---

## 20.10    Point processes and random measures

*This section assumes that the reader is familiar with measure theory. It can however safely be skipped as it is not reused in the rest of the book.*

### 20.10.1    Poisson processes

If $\mathcal{Z}$ is a set, we will denote by $\mathcal{P}_c(\mathcal{Z})$ the set composed with all finite or countable subsets of $\mathcal{Z}$. A point process over $\mathcal{Z}$ is a random variable $S : \Omega \to \mathcal{P}_c(\mathcal{Z})$, i.e., a variable that provides a countable random subset of $\mathcal{Z}$. If $B \subset \mathcal{Z}$ one can then define the counting function $\nu_S(B) = |S \cap B| \in \mathbb{Z} \cup \{+\infty\}$.

A proper definition of such point processes requires some measure theory. Equip $\mathcal{Z}$ with a $\sigma$-algebra $\mathbb{A}$ and consider the set $\mathcal{N}_0$ of integer-valued measures $\mu$ on $(\mathcal{Z}, \mathbb{A})$ such that $\mu(\mathcal{Z}) < \infty$. Let $\mathcal{N}$ be the set formed with all countable sums of measures in $\mathcal{N}_0$. Then a general point process is a mapping $\nu : \Omega \to \mathcal{N}$ such that for all $k \in \mathbb{N} \cup \{+\infty\}$ and all $B \in \mathbb{A}$, the event $\{\nu(B) = k\}$ is measurable. Recall that, for each $B \in \mathbb{A}$, $\nu(B)$ is itself a random variable, that we may denote $\omega \mapsto \nu_\omega(B)$. One then define the intensity of the process as the the function $\mu : B \mapsto E(\nu(B))$.

The following proposition provides an important identity satisfied by such models.

**Theorem 20.16 (Campbell identity)** *Let $\nu$ be a point process with intensity $\mu$. For $\omega \in \Omega$, let $X_\omega : \Omega' \to \mathcal{Z}$ be a random variable with distribution $\nu_\omega$ (defined, if needed, on a different probability space $(\Omega', P')$). Then, for any $\mu$-integrable function $f$:*

$$E(f(X)) = \int_{\mathcal{Z}} f(z) d\mu(z). \qquad (20.22)$$

Here, the expectation of $f(X)$ is over both spaces $\Omega$ and $\Omega'$ and corresponds to the average of $f$. The identity is an immediate consequence of Fubini's theorem.

We will be mainly interested in the family of Poisson point processes. These processes are themselves parametrized by a measure, say $\mu$, on $\mathcal{Z}$ such that $\mu$ is $\sigma$-finite and $\mu(B) = 0$ if $B$ is a singleton. A Poisson process with intensity measure $\mu$ is a point process $\nu$ such that:

(i) If $B_1, \ldots, B_n$ are non-intersecting pairwise, then $\nu(B_1), \ldots, \nu(B_n)$ are mutually independent.

(ii) for all $B$, $\nu(B) \sim Poisson(\mu(B))$.

We take the convention that $\nu(B) = 0$ (resp. $= \infty$) almost surely if $\mu(B) = 0$ (resp. $= \infty$). Note that property (i) also implies that if $g_1, \ldots, g_n$ are measurable functions from $\mathcal{Z}$ to $(0, +\infty)$ such that $g_i g_j = 0$ for $i \neq j$, then the variables $\nu(g_i) = \int_{\mathcal{Z}} g_i(z) d\nu(z)$ are independent.

If $\mu(\mathcal{Z}) < \infty$ (i.e., $\mu$ is finite), one can represent the distribution of a Poisson point process as follows:

$$\nu = \sum_{k=1}^{\nu(\mathcal{Z})} \delta_{X_k}$$

with $\nu(\mathcal{Z}) \sim Poisson(\mu(\mathcal{Z}))$ and, conditional to $\nu(\mathcal{Z}) = N$, $X_1, \ldots, X_N$ are i.i.d. and follow the probability distribution $\bar{\mu} = \mu/\mu(\mathcal{Z})$. This measure can also be identified with the random set $S = \{X_1, \ldots, X_{\nu(\mathcal{Z})}\}$. The assumption that $\mu(\{z\}) = 0$ for any singleton implies that $\nu(\{z\}) = 0$ almost surely. It also ensures that the points $X_1, \ldots, X_N$ are distinct with probability one.

If $\mu$ is $\sigma$-finite, then (by definition), it is a countable sum of finite measures $\mu_1, \mu_2, \ldots$. Then $\nu$ can be generated as the sum of independent $\nu_1, \nu_2, \ldots$, where $\nu_i$ is a Poisson process with intensity $\mu_i$. It can moreover be identified with the countable random set $S = \bigcup_{i=1}^{\infty} S_i$, where $S_i$ is the random set associated with $\nu_i$. Note that, in this construction, one can always assume that the measures $\mu_1, \mu_2, \ldots$ are mutually singular (i.e., $\mu_i(B) > 0$ for some $i$ implies that $\mu_j(B) = 0$ for $j \neq i$).

If we consider a Poisson process on $(0, +\infty) \times \mathcal{Z}$, we can define weighted random measures. Indeed, such a point process takes values in the collection of all sets of the form $\{(w_k, z_k), k \in I\}$ where $I$ is finite or countable. These subsets can be represented as the sum of weighted Dirac masses,

$$\xi = \sum_{k \in I} w_k \delta_{z_k}.$$

To ensure that the points $(z_k, k \in I)$ generated by this process are all different, we need to assume that the intensity $\mu$ of this random process is such that $\mu((0, +\infty) \times \{z\}) = 0$ for all $z \in \mathcal{Z}$. We will refer to $\xi$ as a weighted Poisson process.

In the following, we will consider this class of random measures, with the small addition of allowing for an extra term including a measure supported by a fixed set. More precisely, given a (deterministic) countable subset $\mathcal{I} \subset \mathcal{Z}$, a family of independent random variables $(\rho_z, z \in \mathcal{I})$ and a $\sigma$-finite measure $\mu_o$ such that $\mu_o((0, +\infty) \times \{z\}) = 0$ for all $z \in \mathcal{Z}$, we can define the random measure

$$\xi = \xi_f + \xi_o$$

where $\xi_o$ is a weighted Poisson process with intensity $\mu_o$, assumed independent of $(\rho_z, z \in \mathcal{I})$ and

$$\xi_f = \sum_{z \in \mathcal{I}} \rho_z \delta_z.$$

The subscripts $o$ and $f$ come from the terminology introduced in Kingman [105], which studies "completely random measures," which are a random measures that satisfy point (i) in the definition of a Poisson process. Under mild assumptions, such measures can be decomposed as a sum of a weighted Poisson process (here, $\xi_o$, the ordinary part), of a process with fixed support, (here, $\xi_f$, the fixed part) and of a deterministic measure (which is here taken to be 0).

Let us rapidly check that $\xi$ satisfies property (i). Let $B_1, \ldots, B_n$ be non-overlapping elements of $\mathbb{A}$. Get $g_i(w, z) = w \mathbf{1}_{B_i}(z)$. Then

$$\xi(B_i) = \xi_f(B_i) + \nu_o(g_i)$$

where $\nu_o$ is a Poisson process with intensity $\mu_o$. Since the sets do not overlap, the variables $(\xi_f(B_i), i = 1, \ldots, n)$ are independent, and so are $(\nu_o(g_i), i = 1, \ldots, n)$ since $g_i g_j = 0$ for $i \neq j$. Since $\xi_f$ and $\nu_o$ are, in addition independent, we see that $(\xi(B_i), i = 1, \ldots, n)$ are independent.

The intensity measure of such a process is still defined by

$$\eta(B) = E(\xi(B)) = \sum_{z \in \mathcal{I}} P((\rho_z, z) \in B) + \int_{(0, +\infty) \times B} w \, d\mu_o(w, z)$$

where the last term is an application of Campbell's inequality to the Poisson process $\nu_o$ and the function $g(w, x) = w \mathbf{1}_B(x)$.

### 20.10.2 The gamma process

The main example of such processes in factor analysis is the beta process that will be discussed in the next section. We start, however, with a first example that is closely related with the Dirichlet process, called the gamma process.

In this process, one fixes a finite measure $\pi_0$ on $\mathcal{Z}$ and defines $\mu$ on $(0, +\infty) \times \mathcal{Z}$ by

$$\mu(dw, dz) = cw^{-1}e^{-cw}\pi_0(dz)dw.$$

Because $\mu$ is $\sigma$-finite but not finite (the integral over $t$ diverges at $t = 0$), every realization of $\xi$ is an infinite sum

$$\xi = \sum_{k=1}^{\infty} w_k \delta_{z_k}.$$

The intensity measure of $\xi$ is

$$\eta(B) = c\pi_0(B)\int_0^{+\infty} e^{-cw}dw = \pi_0(B).$$

In particular,

$$\sum_{k=1}^{\infty} w_k = \eta(\mathcal{Z}) = \pi_0(\mathcal{Z}) < \infty.$$

For fixed $B$, the variable $\xi(B)$ follows a Gamma distribution. This can be proved by computing the Laplace transform of $\xi$, $E(e^{-\lambda\xi(B)})$, and identify it to that of a Gamma. To make this computation, consider the point process $\nu_J$ restricted to a interval $J \subset (0, +\infty)$ with $\min(J) > 0$, and $\xi_J$ the corresponding weighted process. Let $m_J(t) = \int_J w^{-1}ce^{-(c+t)w}\,dw$. Then a realization of $\nu_J$ can be obtained by first sampling $N$ from a Poisson distribution with parameter $\mu(J \times \mathcal{Z}) = m_J(0)\pi_0(\mathcal{Z})$ and then sampling $N$ points $(w_i, z_i)$ independently from the distribution $\mu/(m_J(0)\pi_0(\mathcal{Z}))$. This implies that

$$
\begin{aligned}
E(e^{-t\xi_J(B)}) &= \sum_{n=0}^{\infty} e^{-m_J(0)\pi_0(\mathcal{Z})}\frac{(m_J(0)\pi_0(\mathcal{Z}))^n}{n!}\left(\frac{\int_0^\infty e^{-tw\mathbf{1}_B(z)}w^{-1}ce^{-cw}dwd\pi_0(z)}{m_J(0)\pi_0(\mathcal{Z})}\right)^n \\
&= \sum_{n=0}^{\infty} \frac{e^{-m_J(0)\pi_0(\mathcal{Z})}}{n!}\left(\pi_0(B)m_J(t) + (\pi_0(\mathcal{Z}) - \pi_0(B))m_J(0)\right)^n \\
&= e^{\pi_0(B)(m_J(t)-m_J(0))}.
\end{aligned}
$$

Now,

$$m_J(t) - m_J(0) = c\int_J e^{cw}\frac{e^{-tw}-1}{w}dw$$

is finite even when $J = (0, +\infty)$. With a little more work justifying passing to the limit, one finds that, for $J = (0, +\infty)$,

$$E(e^{-t\xi_J(B)}) = \exp\left(\pi_0(B)\int_0^{+\infty} e^{-cw}\frac{e^{-tw}-1}{w}dw\right).$$

Finally, write

$$c\int_0^{+\infty} e^{-cw}\frac{e^{-tw}-1}{w}dw = -c\int_0^{+\infty} e^{-cw}\int_0^t e^{-sw}dsdw$$

$$= -c\int_0^t\int_0^{+\infty} e^{-(s+c)w}dwds$$

$$= -\int_0^t c(s+c)^{-1}ds = -c\log(1+\frac{t}{c}).$$

This shows that

$$E(e^{-t\xi_J(B)}) = \left(1+\frac{t}{c}\right)^{-c\pi_0(B)}$$

which is the Laplace transform of a Gamma distribution with parameters $c\pi_0(B)$ and $c$, i.e., with density proportional to $w^{c\pi_0(B)-1}e^{-cw}$.

As a consequence, the normalized process $\delta = \xi/\xi(\mathcal{Z})$ is a Dirichlet process with intensity $c\pi_0$. Indeed, if $B_1, \ldots, B_n$ is a partition of $\mathcal{Z}$ the family $(\delta(B_1), \ldots, \delta(B_n))$ is the ratio of $n$ independent gamma variables to their sum, which provides a Dirichlet distribution, and this property characterizes Dirichlet processes.

### 20.10.3   The beta process

The definition of the beta process parallels that of the gamma process, with weights taking this time values in $(0, 1)$. Fix again a finite measure $\pi_0$ on $\mathcal{Z}$ and let $\mu_o$ on $(0, +\infty) \times \mathcal{Z}$ be defined by

$$\mu_o(dw, dz) = cw^{-1}(1-w)^{c-1}\pi_0(dz)dw.$$

The associated weighted Poisson process can therefore be represented as a sum

$$\xi_o = \sum_{k=1}^{\infty} w_k\delta_{z_k},$$

and its intensity measure is

$$\eta_o(B) = c\pi_0(B)\int_0^1 (1-t)^{c-1}dw = \pi_0(B).$$

In particular, since $\pi_0$ is finite, we have $\sum_{k=1}^{\infty} w_k < \infty$ almost surely. A beta process is the sum of the process $\xi_o$ and of a fixed set process

$$\xi_f = \sum_{z \in \mathcal{I}} w_z \delta_z$$

where $\mathcal{I}$ is a fixed finite set and $(w_z, z \in \mathcal{I})$ are independent and follow a beta distribution with parameters $(a(z), b(z))$.

If $\mathcal{Z}$ is a space of features, such a process provides a prior distribution on feature selections. It indeed provides, in addition to the deterministic set $\mathcal{I}$, a random countable set $\mathcal{J} \subset \mathcal{Z}$, with a set of random weights $w_z, z \in \mathcal{F} := \mathcal{I} \cup \mathcal{J}$. Given this, one defines the feature process as the selection of a subset $A \subset \mathcal{F}$ where each feature $z$ is selected with probability $w_z$. Because $E(|A|) = \sum_{z \in \mathcal{F}} w_z$ is finite, $A$ is finite with probability 1.

In the same way the Polya urn could be used to sample from a realization of a Dirichlet process without actually sampling the whole process, there exists an algorithm that samples a sequence of feature sets $(A_1, \ldots, A_n)$ from this feature selection process without needing the infinite collection of weights and features associated with a beta process. We assume in the following that the prior process has an empty fixed set. (Non-empty fixed sets will appear in the posterior.)

The first set of features, $A_1$, is obtained as follows according to a Poisson process with intensity $\pi_0$: choose the number $N$ of features in $A_1$ according to a Poisson distribution with parameter $\pi_0(\mathcal{Z})$. Then sample $N$ features independently according to the distribution $\pi_0/\pi_0(\mathcal{Z})$.

Now assume that $n-1$ sets of features $A_1, \ldots, A_n$ have been obtained and we want to sample a new set $A_{n+1}$ conditionally to their observation. Let $\mathcal{J}_n$ be the union of all random features obtained up to this point and $n(z)$, for $z \in \mathcal{J}_n$ the number of times this feature was observed in $A_1, \ldots, A_n$. Then the conditional distribution of the beta process $\xi$ given this observation is still a beta process, with fixed set given by $\mathcal{I} = \mathcal{J}_n$, $(a(z), b(z)) = (n(z), c + n - n(z))$ for $z \in \mathcal{J}_{n-1}$ and base measure $\pi_n = c\pi_0/(c + n)$. This implies that the next set $A_{n+1}$ can be obtained by sampling from the associated feature process. To do this, one first selects features $z \in \mathcal{J}_n$ with probability $n(z)/(c + n)$, then selects additional features $z_1, \ldots, z_N$ independently with distribution $\pi_0/\pi_0(\mathcal{Z})$ where $N$ follows a Poisson distribution with parameter $c\pi_0(\mathcal{Z})/(c + n)$. This is the Indian buffet process, described in Algorithm 20.6 (taking $\pi_0 = \gamma\psi$).

### 20.10.4 Beta Process and feature selection

The beta process can be used as a prior for feature selection within a factor analysis model, as described in the previous paragraph. It is however easier to approximate

it with a model with almost surely finite support. Indeed, letting, for $\epsilon > 0$

$$\mu(dw, dz) = \frac{\Gamma(c+1)}{\Gamma(\epsilon+1)\Gamma(c-\epsilon)} w^{\epsilon-1}(1-w)^{c-\epsilon}\pi_0(dz)dw,$$

one obtains a finite measure since

$$\int_0^{+\infty}\int_{\mathcal{Z}}\mu_o(dw, dz) = \frac{c\gamma}{\epsilon}$$

where $\gamma = \pi_0(\mathcal{Z})$. Note that $\mu$ is normalized so that $E(\xi(B)) = \pi_0(B)$ for $B \subset \mathcal{Z}$.

In this case, the prior generates features by first sampling their number, $p$, randomly according to a Poisson distribution with mean $c\gamma/\epsilon$, then select $p$ probabilities $w_1, \ldots, w_p$ independently using a beta distribution with parameters $\epsilon$ and $c-\epsilon$, and finally attach to each $i$ a feature $z_i$ with distribution $\pi_0/\gamma$. The features associated with a given sample are then obtained by selecting each $z_i$ with probability $w_i$.

We note also that the model described in section 20.9.3 provides an approximation of this prior using a finite number of features. With our notation here, this corresponds to taking $p \gg 1$ and $\epsilon = c\gamma/p$.

# Chapter 21

# Data Visualization and Manifold Learning

## 21.1 Multidimensional scaling

The methods described in this chapter aim at representing a dataset in low dimension, allowing for its visual exploration by summarizing its structure in a user-accessible interface. Unlike factor analysis methods, they do not necessarily attempt at providing a causal model expressing the data as a function of a small number of sources, and generally do not provide a direct mechanism for adding new data to the representation. In addition, all these methods take as input similarity dissimilarity matrices between data points and do not require, say, Euclidean coordinates.

Assuming that a dissimilarity matrix $D = (d_{kl}, k, l = 1, \ldots, N)$ is given, the goals of multidimensional scaling (or MDS) is to determine a small-dimensional Euclidean representation, say $y_1, \ldots, y_N \in \mathbb{R}^p$, such that $\left| y_k - y_l \right|^2 \simeq d_{kl}^2$. We review below two versions of this algorithm, referred to as "similarity" and "dissimilarity" matching.

### 21.1.1 Similarity matching (Euclidean case)

We start with the standard hypotheses of MDS, assuming that the distances $d_{kl}$ derive from a representation in feature space, so that $d_{kl}^2 = \|h_k - h_l\|_H^2$ for some inner-product space $H$ and (possibly unknown) features $h_1, \ldots, h_N$. Note that, since the Euclidean distance is invariant by translation, there is no loss of generality in assuming $h_1 + \cdots + h_N = 0$, which will be done in the following.

We look for a $p$-dimensional representation in the form $y_k = \Phi h_k$ where $\Phi$ is a linear transformation (and we want $y_k$ to be computable directly from dissimilarities, since we do not assume that $h_k$ is known). Since we are only interested in a transformation of the $h_1, \ldots, h_N$, it suffices to compute $\Phi$ in the vector space generated by

them, so that we let

$$\Phi : \mathrm{span}(h_1,\ldots,h_N) \to \mathbb{R}^p,$$

and we want $\Phi$ to (approximately) conserve the norm, i.e., be close to being an isometry.

Because isometries are one-to-one and onto, the existence of an exact isometry would require $V \overset{\triangle}{=} \mathrm{span}(h_1,\ldots,h_N)$ to be $p$-dimensional. The mapping $\Phi$ could then be defined as $\Phi(h) = (\langle h, e_1 \rangle_H, \ldots, \langle h, e_p \rangle_H)$ where $e_1,\ldots,e_p$ is any orthonormal basis of $V$. In the general case, however, where $V$ is not $p$-dimensional or less, one can replace it by a best $p$-dimensional approximation of the training data, leading to a problem similar to PCA in feature space.

Indeed, as we have seen in section 20.1.2, this best approximation can be obtained by diagonalizing the Gram matrix $S$ of $h_1,\ldots,h_N$, which is such that $s_{kl} = \langle h_k, h_l \rangle_H$. (Recall that we assume that $\bar{h} = 0$, so we do not center the data here.) Using the notation in section 20.1.2, let $\alpha^{(1)},\ldots,\alpha^{(p)}$ denote the eigenvectors associated with the $p$ largest eigenvalues, normalized so that $(\alpha^{(i)})^T S \alpha^{(i)} = 1$ for $i = 1,\ldots,p$. One can then take

$$e_i = \sum_{l=1}^{N} \alpha_l^{(i)} h_l$$

and, for $k = 1,\ldots,N$, $j = 1,\ldots,p$:

$$y_k^{(i)} = \lambda_i^2 \alpha_k^{(i)} \tag{21.1}$$

where $\lambda_i^2$ is the eigenvalue associated with $\alpha^{(i)}$.

This does not entirely address the original problem, since the inner products $s_{kl}$ are not given, but only the distances $d_{kl}$, which satisfy

$$d_{kl}^2 = -2s_{kl} + s_{kk} + s_{ll}. \tag{21.2}$$

This provides a linear system of equations in the unknown $s_{kl}$. This system is underdetermined, because $D$ is invariant by any transformation $h_k \mapsto h_k + h_0$ (for a fixed $h_0$), and $S$ is not. However, the assumption $h_1 + \cdots + h_N = 0$ provides the additional constraint needed to provide a unique solution.

Summing (21.2) over $l$, we then get

$$\sum_{l=1}^{N} d_{kl}^2 = N s_{kk} + \sum_{l=1}^{N} s_{ll}. \tag{21.3}$$

Summing this equation over $k$, we find

$$\sum_{k,l=1}^{N} d_{kl}^2 = 2N \sum_{l=1}^{N} s_{ll}.$$

Using this in (21.3), we get

$$s_{kk} = \frac{1}{N} \sum_{l=1}^{N} d_{kl}^2 - \frac{1}{2N^2} \sum_{k,l=1}^{N} d_{kl}^2,$$

and, from (21.2)

$$s_{kl} = -\frac{1}{2}\left( d_{kl}^2 - \frac{1}{N} \sum_{k'=1}^{N} d_{k'l}^2 - \frac{1}{N} \sum_{l'=1}^{N} d_{kl'}^2 + \frac{1}{N^2} \sum_{k',l'=1}^{N} d_{k'l'}^2 \right).$$

If we denote by $D^{\odot 2}$ the matrix formed with the squared distances $d_{kl}^2$, this identity can we rewritten in the simpler form

$$S = -\frac{1}{2}PD^{\odot 2}P \tag{21.4}$$

with $P = \text{Id}_{\mathbb{R}^N} - \mathbb{1}_N \mathbb{1}_N^T / N$.

We now show that this PCA approach to MDS is equivalent to the problem of minimizing

$$F(y) = \sum_{k,l=1}^{N} (y_k^T y_l - s_{kl})^2 \tag{21.5}$$

over all $y_1, \ldots, y_N \in \mathbb{R}^p$ such that $y_1 + \cdots + y_N = 0$, which can be interpreted as matching "similarities" $s_{kl}$ rather than distances. Indeed, letting $\mathcal{Y}$ denote the $N$ by $p$ matrix with rows $y_1^T, \ldots, y_N^T$, we have

$$F(y) = \text{trace}((\mathcal{Y}\mathcal{Y}^T - S)^2).$$

Finding $\mathcal{Y}$ is equivalent to finding a symmetric matrix $M$ of rank $p$ minimizing $\text{trace}((M-S)^2)$. We have, using the trace inequality (theorem 2.1), and letting $\lambda_1^2 \geq \cdots \geq \lambda_N^2$ (resp. $\mu_1^2 \geq \cdots \geq \mu_p^2$) denote the eigenvalues of $S$ (resp. $M$)

$$\text{trace}((M-S)^2) = \text{trace}(M^2) - 2\text{trace}(MS) + \text{trace}(S^2)$$

$$= \sum_{k=1}^{p} \mu_k^4 - 2\text{trace}(MS) + \sum_{k=1}^{N} \lambda_k^4$$

$$\geq \sum_{k=1}^{p} \mu_k^4 - 2\sum_{k=1}^{p} \lambda_k^2 \mu_k^2 + \sum_{k=1}^{N} \lambda_k^2$$

$$= \sum_{k=1}^{p} (\lambda_k^2 - \mu_k^2)^2 + \sum_{k=p+1}^{N} \lambda_k^4$$

$$\geq \sum_{k=p+1}^{N} \lambda_k^4$$

This lower bound is attained when $M$ and $S$ can be diagonalized in the same orthonormal basis with $\lambda_k^2 = \mu_k^2$ for $k = 1,\ldots,p$. So, letting $S = UDU^T$, where $U$ is orthogonal and $D$ is diagonal with decreasing numbers on the diagonal, an optimal $M$ is given by $M = U_p D_p U_p^T$, where $U_p$ is formed with the first $p$ columns of $A$ and $D_p$ is the first $p \times p$ block of $D$. This shows that the matrix $\mathcal{Y} = U_p D^{1/2}$ provides a minimizer of $F$. The matrix $U = [u^{(1)},\ldots,u^{(N)}]$ differs from the matrix $A = [\alpha^{(1)},\ldots,\alpha^{(N)}]$ above through the normalization of its column vectors: we have $S\alpha^{(i)} = \lambda_i^2 \alpha^{(i)}$ with $(\alpha^{(i)})^T S\alpha^{(i)} = 1$ while $Su^{(i)} = \lambda_i^2 u^{(i)}$ with $(u^{(i)})^T Su^{(i)} = \lambda_i^2$ showing that $\alpha^{(i)} = \lambda_i^{-1} u^{(i)}$. This shows that $A_p = U_p D_p^{-1/2}$ so that $\mathcal{Y}$ can also be rewritten as $\mathcal{Y} = A_p D_p$, i.e., $y_k^{(i)} = \lambda_i^2 \alpha_k^{(i)}$, the same expression that was obtained before.

The minimization of $F$ is called *similarity matching*. Clearly, this method can be applied when one starts directly with a matrix of dissimilarities $S$, provided it satisfies $\sum_{l=1}^{N} s_{kl} = 0$ for all $k$. If this is not the case, then interpreting $s_{kl}$ as an inner product $h_k^T h_l$, it is natural to replace $s_{kl}$ by what would give $(h_k - \bar{h})^T (h_l - \bar{h})$, namely, by

$$s'_{kl} = s_{kl} - \frac{1}{N}\sum_{l'=1}^{N} s_{kl'} - \frac{1}{N}\sum_{k'=1}^{N} s_{k'l} + \frac{1}{N^2}\sum_{k',l'=1}^{N} s_{k'l'}.$$

Interestingly, this discussion provides us with yet another interpretation of PCA.

### 21.1.2  Dissimilarity matching

While the minimization of (21.5) did not provide us with a new way of analyzing the data (since it was equivalent to PCA), the direct comparison of dissimilarities, that is, the minimization of

$$G(y) = \sum_{k,l=1}^{N} (|y_k - y_l| - d_{kl})^2$$

over $y_1,\ldots,y_N \in \mathbb{R}^p$, provides a different approach. Since this may be useful in practice and does not bring in much additional difficulty, we will allow for the possibility of weighting the differences in $G$ and consider the minimization of

$$G(y) = \sum_{k,l=1}^{N} w_{kl}(|y_k - y_l| - d_{kl})^2$$

where $W = (w_{kl})$ is a symmetric matrix of non-negative weights. The only additional complexity resulting by adding weights is that the indeterminacy on $y_1,\ldots,y_N$ is that $G(y) = G(y')$ as soon as $y - y'$ is constant on every connected component of the graph associated with the weight matrix $W$, so that the constraint on $y$ should be replaced by

$$\sum_{k \in \Gamma} y_k = 0$$

for any connected component $\Gamma$ of this graph. (If all weights are positive, then the only non-empty connected component is $\{1,\dots,N\}$ and we retrieve our previous constraint $\sum_{k=1}^{N} y_k = 0$.)

Standard nonlinear optimization methods, such as projected gradient descent, may be used to minimize $G$, but the preferred algorithm for MDS uses a stepwise procedure resulting from the addition of an auxiliary variable. Rewrite

$$G(y) = \sum_{k,l=1}^{N} w_{kl}|y_k - y_l|^2 - 2\sum_{k,l=1}^{N} w_{kl}d_{kl}|y_k - y_l| + \sum_{k,l=1}^{N} d_{kl}^2.$$

We have, for $u \in \mathbb{R}^p$:

$$|u| = \max\{z^T u : z \in \mathbb{R}^p, |z| = \mathbf{1}_{u \neq 0}\}.$$

Using this identity, we can introduce auxiliary variables $z_{kl}$, $k,l = 1,\dots,N$ in $\mathbb{R}^p$, with $|z_{kl}| = 1$ if $y_k \neq y_l$ and define

$$\hat{G}(y,z) = \sum_{k,l=1}^{N} w_{kl}|y_k - y_l|^2 - 2\sum_{k,l=1}^{N} w_{kl}d_{kl}(y_k - y_l)^T z_{kl} + \sum_{k,l=1}^{N} d_{kl}^2.$$

We then have

$$G(y) = \min_{z:|z_{kl}|=1 \text{ if } y_k \neq y_k} \hat{G}(y,z).$$

As a consequence, minimizing $G$ in $y$ can be achieved by minimizing $\hat{G}$ in $y$ and $z$ and discarding $z$ when this is done. One can minimize $\hat{G}$ iteratively, alternating minimization in $y$ given $z$ and in $z$ given $y$, both steps being elementary. In order to describe these steps, introduce some matrix notation.

Let $L$ denote the Laplacian matrix of the weighted graph on $\{1,\dots,N\}$ associated with the weight matrix $W$, namely $L = (\ell_{kl}, k,l = 1,\dots,N)$ with $\ell_{kk} = \sum_{k=1}^{N} w_{kl} - w_{kk}$ and $\ell_{kl} = -w_{kl}$ when $k \neq l$. Then,

$$\sum_{k,l=1}^{N} w_{kl}|y_k - y_l|^2 = 2\mathrm{trace}(\mathcal{Y}^T L \mathcal{Y}).$$

Defining $u_k \in \mathbb{R}^p$ by

$$u_k = \sum_{l=1}^{N} w_{kl}d_{kl}(z_{kl} - z_{lk}),$$

and $\mathcal{U} = \begin{pmatrix} u_1^T \\ \vdots \\ u_N^T \end{pmatrix}$, we have

$$\sum_{k,l=1}^{N} w_{kl}d_{kl}(y_k - y_l)^T z_{kl} = \mathrm{trace}(\mathcal{U}^T \mathcal{Y}).$$

With this notation, the optimal matrix $\mathcal{Y}$ must minimize

$$2\mathrm{trace}(\mathcal{Y}^T L \mathcal{Y}) - 2\mathrm{trace}(\mathcal{U}^T \mathcal{Y}).$$

Let $m$ be the number of connected components of the weighted graph. Recall that the matrix $L$ is positive semi-definite and that an orthonormal basis of its null space is provided by vectors, say $e_1,\dots,e_m$, that are constant on each of the $m$ connected components of the graph, so that the constraint on $\mathcal{Y}$ can be written as $e_j^T \mathcal{Y} = 0$ for $j = 1,\dots,m$. Introduce the matrix

$$\hat{L} = L + \sum_{k=1}^{m} e_k e_k^T$$

which is positive definite. Our minimization problem is then equivalent to minimizing

$$2\mathrm{trace}(\mathcal{Y}^T \hat{L} \mathcal{Y}) - 2\mathrm{trace}(\mathcal{U}^T \mathcal{Y}),$$

subject to $e_j^T \mathcal{Y} = 0$ for $j = 1,\dots,m$. The derivative of this function is

$$4\hat{L}\mathcal{Y} - 2\mathcal{U}$$

so that an optimal $\mathcal{Y}$ must satisfy

$$4\hat{L}\mathcal{Y} - 2\mathcal{U} + \sum_{j=1}^{m} e_j \mu_j^T = 0$$

for Lagrange multipliers $\mu_1,\dots,\mu_m \in \mathbb{R}^p$. This shows that

$$\mathcal{Y} = \frac{1}{2}\hat{L}^{-1}\left(\mathcal{U} - \frac{1}{2}\sum_{j=1}^{m} e_j \mu_j^T\right) = \frac{1}{2}\hat{L}^{-1}\mathcal{U} - \frac{1}{4}\sum_{j=1}^{m} e_j \mu_j^T$$

where we have used the fact that $\hat{L}^{-1} e_j = e_j$. We can now identify $\mu_j$ since

$$0 = e_j^T \mathcal{Y} = \frac{1}{2} e_j^T \hat{L}^{-1} \mathcal{U} - \frac{1}{4}\sum_{j'=1}^{m} e_j^T e_{j'} \mu_j^T = \frac{1}{2} e_j^T \mathcal{U} - \frac{1}{4}\mu_j^T$$

so that $\mu_j^T = 2 e_j^T \mathcal{U}$ and the optimal $\mathcal{Y}$ is

$$\mathcal{Y} = \frac{1}{2}\hat{L}^{-1}\mathcal{U} - \frac{1}{2}\sum_{j=1}^{m} e_j e_j^T \mathcal{U}.$$

Note that this expression can be rewritten as

$$\mathcal{Y} = \frac{1}{2} P_L \hat{L}^{-1} \mathcal{U}$$

where $P_L = \mathrm{Id}_{\mathbb{R}^N} - \sum_{k=1}^{N} e_j e_j^T$ is the projection onto the space perpendicular to the null space of $L$ (i.e., the range of $L$). In the case where the graph has a single connected component, one has $m = 1$ and $e_1 = \mathbb{1}_N / \sqrt{N}$ yielding

$$P_L = \mathrm{Id}_{\mathbb{R}^N} - \frac{1}{N} \mathbb{1}_N \mathbb{1}_N^T.$$

The minimization in $z$ given $y$ is straightforward: if $y_k \neq y_k$, then $z_{kl} = (y_k - y_l)/|y_k - y_l|$. If $y_k = y_l$, then one can take any value for $z_{kl}$ and the simplest if of course $z_{kl} = 0$. Using the previous computation, we can summarize a training algorithm for multi-dimensional scaling, called SMACOF for "Scaling by Maximizing a Convex Function" (see, e.g., Borg and Groenen [36] for more details and references).

---

**Algorithm 21.1 (SMACOF)**

Assume that a symmetric matrix of dissimilarities $(d_{kl}, k, l = 1, \ldots, N)$ is given, together with a matrix of weights $(w_{kl}, k, l = 1, \ldots, N)$. Fix a target dimension, $p$. Fix a tolerance constant $\epsilon$.

1. Compute the Laplacian matrix $L$ of the graph associated with the weights, the projection matrix $P_L$ onto the range of $L$ and the matrix $M = (L + \mathrm{Id}_{\mathbb{R}^N} - P_L)^{-1}$.

2. Initialize the algorithm with some family $y_1, \ldots, y_N \in \mathbb{R}^p$ and let $\mathcal{Y} = \begin{pmatrix} y_1^T \\ \vdots \\ y_N^T \end{pmatrix}$.

3. At a given step of the algorithm, let $\mathcal{Y}$ be the current solution and compute, for $k = 1, \ldots, N$:

$$u_k = 2 \sum_{l=1}^{N} w_{kl} d_{kl} \frac{y_k - y_l}{|y_k - y_l|} \mathbb{1}_{y_k \neq y_l}$$

to form the matrix $\mathcal{U} = \begin{pmatrix} u_1^T \\ \vdots \\ u_N^T \end{pmatrix}$.

4. Compute $\mathcal{Y}' = \frac{1}{2} P_L M \mathcal{U}$.

5. If $|\mathcal{Y} - \mathcal{Y}'| \leq \epsilon$, exit and return $\mathcal{Y}'$.

6. Return to step 3.

---

Figure 21.1:    Left:  Multidimensional scaling applied to a 3D curve embedded in a 10-dimensional space retrieves the Euclidean structure. Right: Isomap, in contrasts, identifies the one-dimensional nature of the data.

## 21.2    Manifold learning

The goal of MDS is to map a full matrix of distances into a low-dimensional Euclidean space.  Such a representation, however, cannot address the possibility that the data is supported by a low-dimensional, albeit nonlinear, space.  For example, people leaving on Earth live, for all purposes, on a two-dimensional structure (a sphere), but any faithful Euclidean representation of the world population needs to use the three spatial dimensions.  One may also argue that the relevant distance between points on Earth is not the Euclidean one either (because one would never travel *through* Earth to go from one place to another), but the distance associated to the shortest path on the sphere, which is measured along great circles.

To take another example, the left panel in fig. 21.1 provides the result of applying MDS to a ten-dimensional dataset obtained by applying a random ten-dimensional rotation to a curve supported by a three-dimensional torus.  MDS indeed retrieves the correct curve structure in space, which is three dimensional.  However, for a person "living" on the curve, the data is one-dimensional, a fact that is captured by the Isomap method that we now describe.

### 21.2.1    Isomap

Let us return to the example of people living on the spherical Earth.  One can define the distance between two points on Earth either as the shortest length a person would have to travel (say, by plane) to go from one point to the other (that we can call the *intrinsic distance*), or simply the *chordal distance* in 3D space between the two points.  The first one is obviously the most relevant to the spherical structure of the Earth, but the second one is easier to compute given the locations of the points in

space.

For typical datasets, the geometric structure of the data (e.g., that it is supported by a sphere) is unknown, and the only information that is available is their chordal distance in an ambient space (which can be very large). An important remark, however is that, when the points are close to each other, the two distances can be expected to be similar, if we assume that the geometry of the set supporting the data is locally linear (e.g., that it is, like the sphere, a "submanifold" of the ambient space, with small neighborhoods of any data point well approximated, at first order, by points on a tangent space). Isomap uses this property, only trusting small distances in the matrix $D$, and infers large distances by adding the costs resulting from traveling from data points to nearby data points.

Fix an integer $c$. Given $D$, the $c$-nearest neighbor graph on $\mathcal{V} = \{1,\ldots,N\}$ places an edge between $k$ and $l$ if and only if $d_{k,l}$ is among the $c$ smallest values in $\{d_{kl'}, l' \neq k\}$ neighbors or $x_l$ among the $c$ smallest values in $\{d_{k'l}, k' \neq l\}$. We will write $k \sim_c l$ to indicate that there exists an edge between $k$ and $l$ in this graph. One then defines the geodesic distance on the graph as

$$d_{kl}^{(*)} = \min\left\{\sum_{j=1}^m d_{k_{j-1}k_j} : k_0,\ldots,k_m \in \{1,\ldots,N\}, k_0 = k \sim_c k_1 \sim_c \cdots \sim_c k_{m-1} \sim_c k_m = l, m \geq 0\right\}.$$

This geodesic distance can be computed incrementally as follows. First define $d_{kl}^{(1)} = |x_k - x_l|$ if $k \sim_c l$ and $d_{kl}^{(1)} = +\infty$ otherwise (and also let $d_{kk}^{(1)} = 0$). Then, given $d^{(n-1)}$, define

$$d_{kl}^{(n)} = \min\left\{d_{kl'}^{(n-1)} + d_{ll'}^{(1)} \; l' = 1,\ldots,N\right\}$$

until the entries stabilize, i.e., $d^{(n+1)} = d^{(n)}$, in which case one has $d^{(*)} = d^{(n)}$. The validity of the statement can be easily proved by checking that

$$d_{kl}^{(n)} = \min\left\{\sum_{j=1}^n d_{k_{j-1}k_j}^{(1)} : k_0,\ldots,k_n \in \{1,\ldots,N\}, k_0 = k, k_n = l\right\},$$

which can be done by induction, the details being left to the reader. It should also be clear that the procedure will stabilize after no more than $N$ steps.

Once the distance is computed, Isomap then applies standard MDS, resulting in a straightened representation of the data like in fig. 21.1. Another example is provided in fig. 21.2, where, this time, the input curve is closed and cannot therefore be represented as a one-dimensional structure. One can note, however, that, even in this case, Isomap still provides some simplification of the initial shape of the data.
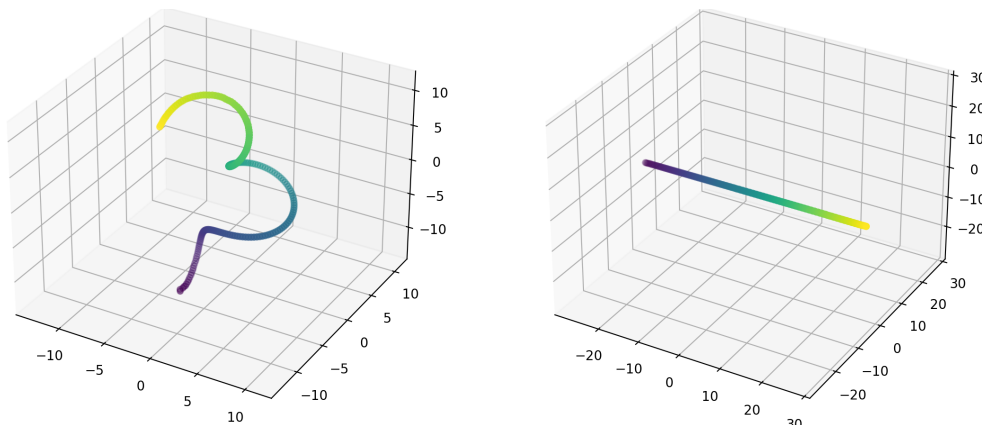
Figure 21.2: Left: Multidimensional scaling applied to a 3D curve embedded in a 10-dimensional space retrieves the Euclidean structure. Right: Isomap, in contrasts, identifies the one-dimensional nature of the data.

### 21.2.2  Local Linear Embedding

Local linear embedding (LLE) exploits in a different way the fact that manifolds are locally well approximated by linear spaces. Like Isomap, it starts also with building a $c$-nearest-neighbor graph on $\{1,\dots,k\}$. Assume, for the sake of the discussion, that the distance matrix is computed for possibly unobserved data $T = (x_1,\dots,x_N)$. Letting $\mathcal{N}_k$ denote the indices of the nearest neighbors of $k$ (excluding $k$ itself), the basic assumption is that $x_k$ should approximately lie in the affine space generated by $x_l, l \in \mathcal{N}_k$. Expressed in barycentric coordinates, this space is defined by

$$\mathcal{T}_k = \left\{ \sum_{l \in \mathcal{N}_k} \rho^{(l)} x_l : \rho \in \mathbb{R}^{\mathcal{N}_k}, \sum_{l \in \mathcal{N}_k} \rho^{(l)} = 1 \right\},$$

and $\mathcal{T}_k$ can be interpreted as an approximation of the tangent space at $x_k$ to the data manifold. Optimal coefficients $(\rho^{()}kl, k = 1,\dots,N, l \in \mathcal{N}_k)$ providing the representation of $x_k$ in that space can be estimated by minimizing, for all $k$

$$\left| x_k - \sum_{l \in \mathcal{N}_k} \rho_k^{(l)} x_l \right|^2$$

subject to $\sum_{l \in \mathcal{N}_k} \rho_k^{(l)} = 1$. This is a simple least-square program. Let $c_k = |\mathcal{N}_k|$ ($c_k = c$ in the absence of ties). Order the elements of $\mathcal{N}_k$ to represent $\rho_k^{(l)}, l \in \mathcal{N}_k$ as a vector denoted $\boldsymbol{\rho}_k \in \mathbb{R}^{c_k}$. Similarly, let $S_k$ be the Gram matrix associated with $x_l, l \in \mathcal{N}_k$ formed with all inner products $x_{l'}^T x_l, l, l' = 1,\dots,N$ and let $\boldsymbol{r}_k$ be the vector composed with products $x_k^T x_l, l \in \mathcal{N}_k$. Assume that $S_k$ is invertible, which is generally true if

$c < d$, unless the neighbors are exactly linearly aligned. Then, the optimal $\rho_k$ and the Lagrange multiplier $\lambda$ for the constraint are given by

$$\begin{pmatrix} \rho_k \\ \lambda \end{pmatrix} = \begin{pmatrix} S_k & \mathbb{1}_{c_k} \\ \mathbb{1}_{c_k}^T & 0 \end{pmatrix}^{-1} \begin{pmatrix} r_k \\ 1 \end{pmatrix}. \tag{21.6}$$

If $S_k$ is not invertible, the problem is under-constrained and one of its solutions can be obtained by replacing the inverse above by a pseudo-inverse.

The low-dimensional representation of the data, still denoted $(y_1, \ldots, y_N)$ with $y_k \in \mathbb{R}^p$ is then estimated so that the relative position of $y_k$ to its neighbors is the same as that of $x_k$, i.e., so that

$$y_k \simeq \sum_{l \in \mathcal{N}_k} \rho_k^{(l)} y_l.$$

These vectors are estimated by minimizing

$$F(y) = \sum_{k=1}^{N} \left| y_k - \sum_{l \in \mathcal{N}_k} \rho_k^{(l)} y_l \right|^2.$$

Obviously, some additional constraints are needed to avoid the trivial solution $y_k = 0$ for all $k$. Also, replacing all $y_k$'s by $y_k' = R y_k + b$ where $R$ is an orthogonal transformation in $\mathbb{R}^p$ and $b$ is a translation does not change the value of $F$, so there is no loss of generality in assuming that $\sum_{k=1}^{N} y_k = 0$ and that $\sum_{k=1}^{N} y_k y_k^T = D_0$, a diagonal matrix. However, if one lets $y_k' = D y_k$ where $D$ is diagonal, then

$$F(y) = \sum_{i=1}^{p} D_{ii}^2 \sum_{k=1}^{N} \left( y_k^{(i)} - \sum_{l \in \mathcal{N}_k} \rho_k^{(l)} y_l^{(i)} \right)^2.$$

This shows that one should not allow the diagonal coefficients of $D_0$ to be chosen freely, since otherwise the optimal solution would require to take this coefficient to 0. So $D_0$ should be a fixed matrix, and by symmetry, it is natural to take $D_0 = \mathrm{Id}_{\mathbb{R}^p}$. (Any other solution—for a different $D_0$—can then be obtained by rescaling independently the coordinates of $y_1, \ldots, y_N$.)

Extend $\rho_k^{(l)}$ to an $N$-dimensional vector by taking $\rho_k^{(k)} = -1$ and $\rho_k^{(l)} = 0$ if $l \neq k$ and $l \notin \mathcal{N}_k$. We can write

$$F(y) = \sum_{k=1}^{N} \left| \sum_{l=1}^{N} \rho_k^{(l)} y_l \right|^2.$$

Expanding the square, this is

$$F(y) = \sum_{l,l'=1}^{N} w_{ll'} y_l^T y_{l'}$$

with $w_{ll'} = \sum_{k=1}^{N} \rho_k^{(l)} \rho_k^{(l')}$. Introducing the matrix $\mathcal{W}$ with entries $w_{kl}$ and the $N \times p$ matrix $\mathcal{Y} = \begin{pmatrix} y_1^T \\ \vdots \\ y_N^T \end{pmatrix}$, we have the simple expression

$$F(y) = \text{trace}(\mathcal{Y}^T \mathcal{W} \mathcal{Y}).$$

Note that the constraints are $\mathcal{Y}^T \mathcal{Y} = \text{Id}_{\mathbb{R}^p}$ and $\mathcal{Y}^T \mathbb{1}_N = 0$. Without this last constraint, we know that an optimal solution is provided by $\mathcal{Y} = [e_1, \ldots, e_p]$ where $e_1, \ldots, e_p$ provide an orthonormal family of eigenvectors associated to the $p$ smallest eigenvalues of $\mathcal{W}$ (this is a consequence of corollary 2.4). To handle the additional constraint, it suffices to note that $\mathcal{W} \mathbb{1}_N = 0$, so that $\mathbb{1}_N$ is a zero eigenvector. Given this, it suffices to compute $p + 1$ eigenvectors associated to smallest eigenvalues of $\mathcal{W}$, $e_1, \ldots, e_{p+1}$, with the condition that $e_1 = \pm \mathbb{1}_N / \sqrt{N}$ (which is automatically satisfied unless 0 is a multiple eigenvalue of $\mathcal{W}$) and let

$$\mathcal{Y} = [e_2, \ldots, e_{p+1}].$$

Note that $e_2, \ldots, e_{p+1}$ are also the $p$ smallest eigenvectors of $\mathcal{W} + \lambda \mathbb{1} \mathbb{1}^T$ for any large enough $\lambda$, e.g., $\lambda > \text{trace}(\mathcal{W})/N$.

LLE is summarized in the following algorithm.

---

**Algorithm 21.2 (Local linear embedding)**
The input of the algorithm is

(i) Either a training set $T = (x_1, \ldots, x_N)$, or its Gram matrix $S$ containing all inner products $x_k^T x_l$ (or more generally inner products in feature space), or a dissimilarity matrix $D = (d_{kl})$.

(ii) An integer $c$ for the graph construction.

(iii) An integer $p$ for the target dimension.

(1) If not provided in input, compute the Gram matrix $S$ and distance matrix $D$ (using (21.2) and (21.4)).

(2) Build the $c$-nearest-neighbor graph associated with the distances. Let $\mathcal{N}_k$ be the set of neighbors of $k$, with $c_k = |\mathcal{N}_k|$.

(3) For $k = 1, \ldots, N$, let $S_k$ be the sub-matrix of $S$ matrix associated with $x_l, l \in \mathcal{N}_k$ and compute coefficients $\rho_k^{(l)}, l \in \mathcal{N}_k$ stacked in a vector $\rho_k \in \mathbb{R}^{c_k}$ by solving (21.6).

(4) Form the matrix $\mathcal{W}$ with entries $w_{ll'} = \sum_{k=1}^{N} \rho_k^{(l)} \rho_k^{(l')}$ with $\rho$ extended so that $\rho_k^{(k)} = -1$ and $\rho_k^{(l)} = 0$ if $l \neq k$ and $l \notin \mathcal{N}_k$.

Figure 21.3: Local linear embedding with target dimension 3 applied to the data in fig. 21.1 and fig. 21.2.

(5) Compute the first $p + 1$ eigenvectors, $e_1, \ldots, e_{p+1}$, of $\mathcal{W}$ (associated with smallest eigenvalues) arranging for $e_1$ to be proportional to $\mathbb{1}_N$.

(6) Set $y_k^{(i)} = e_{i+1}^{(k)}$ for $i = 1, \ldots, p$ and $k = 1, \ldots, N$.

---

The results of LLE applied to the datasets described in fig. 21.1 and fig. 21.2 are provided in fig. 21.3.

**Remark 21.1** We note that, for both Isomap and LLE, the $c$-nearest-neighbors graph can be replaced by the graph formed with edges between all pairs of points that are at distance less than $\epsilon$ from each other, for a chosen $\epsilon > 0$, with no change in the algorithms.

These parameters ($c$ or $\epsilon$) must be chosen carefully and may have an important impact on the output of the algorithm. Choosing them too small would not allow for a correct estimation of distances in Isomap (with possibly some of them being infinite if the graph has more than one connected component), or of the linear approximations in LLE. However, choosing them too large may break the basic hypothesis that the data is locally Euclidean or linear that form the basic principles of these algorithms. ◆

### 21.2.3 Graph Embedding

Both Isomap and LLE are based on the construction of a nearest-neighbor graph based on dissimilarity data and the conservation of some of its geometric features when deriving a small-dimensional representation. For LLE, a weight matrix $\mathcal{W}$ was first estimated based on optimal linear approximations of $x_k$ by its neighbors, and

the representation was computed by estimating the eigenvectors associated with the smallest eigenvalues of $\mathcal{W}$ (excluding the eigenvector proportional to $\mathbb{1}$). However, both methods were motivated by the intuition that the dataset was supported by a continuous small-dimensional manifold. We now discuss methods that are solely motivated by the discrete geometry of a graph, for which we use tools that are similar to our discussion of graph clustering in section 19.5.

Adapting the notation in that section to the present one, we start with a graph with $N$ vertices and weights $\beta_{kl}$ between these vertices (such that $\beta_{ll} = 0$) and we form the Laplacian operator defined by, for any vector $u \in \mathbb{R}^N$:

$$\frac{1}{2}\|u\|_{H^1} = \frac{1}{2}\sum_{l,l'=1}^{N}\beta_{ll'}(u^{(l)} - u^{(l')})^2 = u^T L u,$$

so that $L$ is identified as the matrix with coefficients $\ell_{ll'} = -\beta_{ll'}$ for $l \neq l'$ and $\ell_{ll} = \sum_{l'=1}^{N}\beta_{ll'}$. The matrix $\mathcal{W}$ that was obtained for LLE coincides with this graph Laplacian if one lets $\beta_{ll'} = -w_{ll'}$ for $l \neq l'$, since we have $\sum_{l'=1}^{N}w_{ll'} = 0$. The usual requirement that weights are non-negative is no real loss of generality, because in LLE (and in the Graph embedding method above), one is only interested in eigenvectors of $\mathcal{W}$ (or $L$ below) that are perpendicular to $\mathbb{1}$, and those remain the same if one replaces $\mathcal{W}$ by

$$\mathcal{W} - a\mathbb{1}_N\mathbb{1}_N^T + Na\,\mathrm{Id}_{\mathbb{R}^N}$$

which has negative off-diagonal coefficients $\tilde{w}_{ll'} = w_{ll'} - a$ for large enough $a$.

In graph (or Laplacian) embedding, the starting point is a weighted graph on $\{1,\ldots,N\}$ with edge weights $\beta_{ll'}$ interpreted as similarities between vertexes. These weights may or may not be deduced from measures of dissimilarity $(d_{ll'}, k, l = 1,\ldots,N)$ which themselves may or may not be computed as distances between training data $x_1,\ldots,x_N$. If one starts with dissimilarities, it is typical to use simple transformations to compute edge weights, and one the most commonly used is

$$\beta_{ll'} = \exp(-d_{ll'}^2/2\tau^2)$$

for some constant $\tau$. These weights are usually truncated, replacing small values by zeros (or the computation is restricted to nearest neighbors), to ensure that the resulting graph is sparse, which speeds up the computation of eigenvectors for large datasets.

Given a target dimension $p$, the graph is then represented as a collection of points $y_1,\ldots,y_N \in \mathbb{R}^p$, where $y_k$ is associated to vertex $k$. For this purpose, one needs to compute the first $p+1$ eigenvectors, $e_1,\ldots,e_{p+1}$, of the graph Laplacian, with the requirement that $e_1 = \pm\mathbb{1}_N/\sqrt{N}$. (This is always possible and can be achieved numerically by computing eigenvectors of $L + c\mathbb{1}\mathbb{1}^T$ for large enough $c$.) The graph representation

is then given by $y_k^{(l)}i = e_{i+1}^{(k)}$ for $i = 1, \ldots, p$ and $k = 1, \ldots, N$. Note that these are exactly the same operations as those described in steps 4 and 5 of the LLE algorithm.

One way to interpret this construction is that $e_2, \ldots, e_{p+1}$ (the coordinate functions for the representation $y_1, \ldots, y_N$) minimize

$$\sum_{j=1}^{p} \|e_i\|_{H_1}^2$$

subject to $e_2, \ldots, e_{p+1}$ being perpendicular to each other and perpendicular to the constant functions (these constraints being justified for the same reasons as those discussed for LLE). Small $H^1$ semi-norms being associated with smoothness on the graph, we see that we are looking for the smoothest zero-mean representation of the data.

Based on our discussion of LLE, we can make an alternative interpretation by introducing a symmetric square root $R$ of the Laplacian matrix $L$ or any matrix such that $RR^T = L$. Writing $R = [\rho_1, \ldots, \rho_N]$, one has

$$L = \sum_{k=1}^{N} \rho_k \rho_k^T$$

and $\sum_{k=1}^{N} \rho_k = 0$. With this notation, we can interpret Laplacian embedding as the minimization of

$$\sum_{k=1}^{N} \left| \sum_{l=1}^{N} \rho_k^{(l)} y_l \right|^2$$

(subject to previous orthogonality constraints). In other terms, $y_1, \ldots, y_N$ are determined so that the linear relationships

$$\rho_{kk} y_k = -\sum_{l \neq k} \rho_k^{(l)} y_l$$

are satisfied, which is similar to the LLE condition, without the requirement that $\rho_k(k) = 1$.

An alternate requirement that could have been made for LLE is that $\sum_{l=1}^{N} (\rho_k^{(l)})^2 = 1$ for all $k$. Instead of having to solve a linear system in step 2 of Algorithm 21.2, one would then compute an eigenvector with smallest eigenvalue of $S_k$. For graph embedding, this constraint can be enforced by modifying the Laplacian matrix, since $\sum_{l=1}^{N} (\rho_k^{(l)})^2$ is just the $(k, k)$ coefficient of $RR^T$. Given this, let $D$ be the diagonal matrix formed by the diagonal elements of $L$, and define the so-called "symmetric Laplacian" $\tilde{L} = D^{-1/2} L D^{-1/2}$. One obtain an alternative, and popular, graph embedding method by replacing $e_1, \ldots, e_{p+1}$ above by the first $p$ eigenvectors of $\tilde{L}$.

Another interpretation of this representation can be based on the random walk associated with the graph structure. Consider the random process $t \mapsto q(t)$ defined as follows. The initial position, $q(0)$ is selected according to some arbitrary distribution, say $\pi_0$. Conditional to $q(t) = k$, the next position is determined by setting random waiting times $\tau_{kl}$, each distributed as an exponential distribution with rate $\beta_{kl}$ (or expectation $1/\beta_{kl}$), and the process moves to the position $l$ for which $\tau_{kl}$ is smallest after waiting for that time. Let $P(t)$ be the matrix with coefficients $P(t,k,l) = P(q(t+s) = l \mid q(s) = k)$. Then, one has

$$P(t) = e^{-tL}$$

where the right-hand side is the matrix exponential. If $\lambda_1 = 0 \le \lambda_2 \le \cdots \le \lambda_N$ are the eigenvalues of $L$ with corresponding eigenvectors $e_1, \ldots, e_N$, then

$$P(t) = \sum_{i=1}^{N} e^{-t\lambda_i} e_i e_i^T$$

In particular, restricting the first eigenvectors of $L$ provides an approximation of this stochastic process, i.e.,

$$P(t) \simeq \frac{\mathbb{1}\mathbb{1}^T}{N} + \sum_{i=1}^{p} e^{-t\lambda_{i+1}} y(i) y(i)^T.$$

We could also have considered the discrete-time version of the walk, for which, considering integer times $t \in \mathbb{N}$,

$$P(q(t+1) = l \mid q(t) = k) = \begin{cases} \dfrac{\beta_{kl}}{\sum_{l'=1, l' \neq k}^{N} \beta_{kl'}} & \text{if } l \neq k \\ 0 & \text{if } l = k \end{cases}$$

Introducing the matrix $B$ of similarities $\beta_{kl}$ (with zero on the diagonal) and the diagonal matrix $D$ with coefficients $d_{kk} = \sum_{l=1, l \neq k}^{N} \beta_{kl}$, the r.h.s. of the previous equation is the $k, l$ entry of the matrix $\tilde{P} = D^{-1}B$. Then, for any integer $s$, $P(q(t+s) = l \mid q(t) = k)$ is the $k, l$ entry if $\tilde{P}^s = D^{-1/2}(D^{-1/2}BD^{-1/2})^s D^{1/2}$.

The Laplacian matrix $L$ is given by $L = D - B$. The normalized Laplacian is

$$\bar{L} = D^{-1/2}LD^{-1/2} = \text{Id}_{\mathbb{R}^N} - D^{-1/2}BD^{-1/2}$$

so that

$$\tilde{P}^s = D^{-1/2}(\text{Id}_{\mathbb{R}^N} - \bar{L})^s D^{1/2}.$$

If one introduces the eigenvectors $\bar{e}_1, \ldots, \bar{e}_N$ of the normalized Laplacian, still associated with non-decreasing eigenvalues $\bar{\lambda}_1 = 0, \ldots, \bar{\lambda}_N$, and arranges without loss of generality that $\bar{e}_1 \propto D^{1/2} \mathbb{1}_N$, then

$$\tilde{P}^s = D^{-1/2} \left( \sum_{i=1}^{N} (1 - \bar{\lambda}_i)^s \bar{e}_i \bar{e}_i^T \right) D^{1/2}.$$

This shows that, for $s$ large enough, the transitions of this Markov chain are well approximated by its first terms, suggesting using the alternative representation based on the normalized Laplacian:

$$\bar{y}_k(i) = \bar{e}_{i+1}(k).$$

Both representations (using normalized or un-normalized Laplacians) are commonly used in practice.

### 21.2.4 Stochastic neighbor embedding

**General algorithm**

Stochastic neighbor embedding (SNE, Hinton and Roweis [90]), and its variant (t-SNE, Maaten and Hinton [122]) have become a popular tool for the visualization of high-dimensional data based on dissimilarity matrices. One of the key contributions of this algorithm is to introduce a local data rescaling step, that allows for visualization of more homogeneous point clouds.

Assume that dissimilarities $D = (d_{kl}, k, l = 1, \ldots, N)$ are observed. The basic principle in SNE is to deduce from the dissimilarities a family of $N$ probability distributions on $\{1, \ldots, N\}$, that we will denote $\pi_k$, $k = 1, \ldots, N$, with the property that $\pi_k(k) = 0$. The computation of these probabilities include the local normalization step, and we will return to this later. Given the $\pi_k$'s, one then estimate low-dimensional representations $\boldsymbol{y} = (y_1, \ldots, y_N)$ such that $\pi_k \simeq \psi_k$ where $\psi_k$ is given by

$$\psi_k(l; \boldsymbol{y}) = \frac{\exp\left(-\beta\left(|y_k - y_l|^2\right)\right)}{\sum_{l'=1, l' \neq k}^{N} \exp\left(-\beta\left(|y_k - y_{l'}|^2\right)\right)} \mathbf{1}_{l \neq k}.$$

Here, $\beta : [0, +\infty) \to [0, +\infty)$ is an increasing differentiable function that tends to $+\infty$ at infinity. The derivative is denoted $\partial\beta$. The original version of SNE [90] uses $\beta(t) = t$ and t-SNE [122] takes $\beta(t) = \log(1 + t)$.

The determination of the representation can then be performed by minimizing a measure of discrepancy between the probabilities $\pi_k$ and $\psi_k$. In Hinton and Roweis [90], it is suggested to minimize the sum of Kullback-Liebler divergences, namely

$$\sum_{k=1}^{N} KL(\pi_k \| \psi_k(\cdot; \boldsymbol{y}))$$

or, equivalently, to maximize

$$F(\mathbf{y}) = \sum_{k,l=1}^{N} \pi_k(l) \log \psi_k(l; \mathbf{y})$$

$$= -\sum_{k,l=1}^{N} \beta(|y_k - y_l|^2)\pi_k(l) + \sum_{k=1}^{N} \log\left(\sum_{l=1,l\neq k}^{N} \exp(-\beta(|y_k - y_l|^2))\right)$$

The gradient of this function can be computed by evaluating the derivative at $\epsilon = 0$ of $f : \epsilon \mapsto F(\mathbf{y} + \epsilon \mathbf{h})$. This computation gives

$$f'(0) = -2\sum_{k,l=1}^{N} \partial\beta(|y_k - y_l|^2)(y_k - y_l)^T(h_k - h_l)\pi_k(l)$$

$$+ 2\sum_{k=1}^{N}\sum_{l=1}^{N} \partial\beta(|y_k - y_l|^2)(y_k - y_l)^T(h_k - h_l)\psi_k(l;\mathbf{y})$$

$$= -2\sum_{k=1}^{N} h_k^T \sum_{l=1}^{N} \partial\beta(|y_k - y_l|^2)(y_k - y_l)(\pi_k(l) + \pi_l(k) - \psi_k(l;\mathbf{y}) - \psi_l(k;\mathbf{y}))$$

This shows that

$$\partial_{y_k} F(\mathbf{y}) = -2\sum_{l=1}^{N} \beta(|y_k - y_l|^2)(y_k - y_l)(\pi_k(l) + \pi_l(k) - \psi_k(l;\mathbf{y}) - \psi_l(k;\mathbf{y})).$$

This is a rather simple expression that can be used with any first-order optimization algorithm to maximize $F$. The algorithm in Hinton and Roweis [90] uses gradient ascent with momentum, namely iterating

$$\mathbf{y}^{(n+1)} = \mathbf{y}^{(n)} + \gamma \nabla F(\mathbf{y}^{(n)}) + \alpha^{(n)}(y^{(n)} - y^{(n-1)})$$

Choosing $\alpha^{(n)} = 0$ provides standard gradient ascent with fixed gain $\gamma$ (of course, other optimization methods may be used). The momentum can be interpreted, in a loose sense, as a "friction term".

A variant of the algorithm replaces the node-dependent probabilities $\pi_k$ by a single, symmetric, joint distribution $\bar\pi$ on $\{1,\ldots,N\}^2$, $(k,l) \mapsto \bar\pi(k,l)$, satisfying $\bar\pi(k,k) = 0$ and $\bar\pi(k,l) = \bar\pi(l,k)$. The target distribution $\bar\psi$ then becomes

$$\bar\psi(k,l;\mathbf{y}) = \frac{\exp(-\beta(|y_k - y_l|^2))}{\sum_{k'\neq l'=1}^{N} \exp(-\beta(|y_{k'} - y_{l'}|^2))}.$$

With such a choice, the objective function has a simpler form, namely minimizing $KL(\bar{\pi}\|\bar{\psi}(\cdot,y))$ or maximizing the expected likelihood

$$\bar{F}(y) = \sum_{k,l=1}^{N} \bar{\pi}(k,l)\log\bar{\psi}(k,l;y) = -\sum_{k,l=1}^{N}\beta(|y_k-y_l|^2)\bar{\pi}(k,l) + \log\left(\sum_{k\neq l=1}^{N}\exp(-\beta(|y_k-y_l|^2))\right).$$

The gradient of this symmetric version of $F$ can be computed similarly to the previous one and is given by

$$\partial_{y_k}\bar{F}(y) = -4\sum_{l=1}^{N}\partial\beta(|y_k-y_l|^2)(y_k-y_l)(\bar{\pi}(k,l)-\bar{\psi}(k,l;y)).$$

**Setting initial probabilities**

The probabilities $\pi_k(l)$ or $\bar{\pi}(k,l)$ are deduced from the dissimilarities as

$$\pi_k(l) = \frac{e^{-d_{kl}^2/2\sigma_k^2}}{\sum_{l'=1,l'\neq k}^{N}e^{-d_{kl'}^2/2\sigma_k^2}}$$

for $l\neq k$ and

$$\bar{\pi}(k,l) = \frac{\pi_k(l)+\pi_l(k)}{2n}.$$

The coefficients $\sigma_k^2$, $k=1,\dots,N$ operate the local normalization, justifying, in particular, the parameter-free expression chosen for $\psi$ and $\bar{\psi}$. These coefficients are estimated so as to adjust the entropies of all $\pi_k$ to a fixed value, which is a parameter of the algorithm. Note that, letting $t=1/2\sigma_k^2$ and $H(\pi_k) = -\sum_{l=1}^{N}\pi_k(l)\log\pi_k(l)$,

$$\partial_t H(\pi_k) = -\sum_{l=1}^{N}\partial_t\pi_k(l)\log\pi_k(l) - \sum_{l=1}^{N}\partial_t\pi_k(l)$$

$$= -\sum_{l=1}^{N}\partial_t\pi_k(l)\log\pi_k(l)$$

Now

$$\partial_t\log\pi_k(l) = -d_{kl}^2 + \bar{d}_k^2$$

with $\bar{d}_k^2 = \sum_{l'=1}^{N}d_{kl'}^2\pi_k(l')$. Writing $\partial_t\pi_k(l) = \pi_k(l)\partial_t\log\pi_k(l)$, we have

$$\partial_t H(\pi_k) = \sum_{l=1}^{N}(d_{kl}\log\pi_k(l))\pi_k(l) - \bar{d}_k\sum_{l=1}^{N}\pi_k(l)\log\pi_k(l).$$

Using Schwartz inequality, we see that $\partial_t H(\pi_k) \leq 0$ so that $H(\pi_k)$ is decreasing as a function of $t$, i.e., increasing as a function of $\sigma_k^2$. When $\sigma_k^2 \to 0$, $\pi_k$ converges to the uniform distribution on the set of nearest neighbors of $k$ (the indexes $l \neq k$ such that $d_{kl}^2$ is minimal) and, letting $\nu_k$ denote their number, which is typically equal to 1, $H(\pi_k)$ converges to $\log \nu_k$. When $\sigma_k^2$ tends to infinity, $\pi_k$ converges to the uniform distribution over indexes $l \neq k$, whose entropy is $\log(N-1)$. This shows that $e^{H(\pi_k)}$, which is called the *perplexity* of $\pi_k$ can take any value between $\nu_k$ and $N-1$. The common target value of the perplexity can therefore be taken anywhere between $\max_k \nu_k$ and $N-1$. In Maaten and Hinton [122], it is recommended to choose a value between 5 and 50.

**Remark 21.2** The complexity of the computation of the gradient of the objective function (either $F$ or $\bar{F}$ ) scales like the square of the size of the training set, which may be prohibitive when $N$ is large. In Van Der Maaten [193], an accelerated procedure, that involves an approximation of the gradient is proposed. (This procedure is however limited to representations in dimensions 2 or 3.)                                 ♦

### 21.2.5   Uniform manifold approximation and projection (UMAP)

UMAP is similar in spirit to t-SNE, with a few important differences that result in a simpler optimization problem and faster algorithms. Like Isomap, the approach is based on matching distances between the high-dimensional data and the low-dimensional representation. But while Isomap estimates a unique distance on the whole training set (the geodesic distance on the nearest-neighbor graph), UMAP estimates as many "local distances" as observations before "patching" them to form the final representation.

The goal of transporting possibly non-homogeneous locally defined objects on initial data to a homogeneous low-dimensional visualization is what makes UMAP similar to t-SNE. The difference is that t-SNE transports local probability distributions, while UMAP transports metric spaces. More precisely, given distances $(d_{kl}, k, l = 1, \ldots, N)$ and an integer $m$ provided as input, the algorithm builds, for each $k = 1, \ldots, N$ a (pseudo-)metric $\delta_k$ on the associated data graph by letting

$$\delta^{(k)}(k,l) = \delta^{(k)}(l,k) = \frac{1}{\sigma_k}\left(d_{kl} - \min_{l' \neq k} d_{kl'}\right)$$

if $l$ is among the $m$ nearest neighbors of $k$, where $m$ is a parameter of the algorithm, with all other distances being infinite. The normalization parameter $\sigma_k$ has a role similar to that of the same parameter in t-SNE in that it tends to make the representation homogeneous. Here, it is computed such that

$$\sum_l \exp(-\delta^{(k)}(l,l')) = \log_2 m.$$

Each such metric provides a weighted graph structure on $\{1,\ldots,N\}$ by defining weights $w_{ll'}^{(k)} = \exp(-\delta^{(k)}(l,l'))$. In UMAP, these weights are interpreted in the framework of *fuzzy sets*, where a fuzzy set is defined by a pair $(A,\mu)$ where $A$ is a set and $\mu$ a function $\mu : A \to [0,1]$ [210]. The function $\mu$ is called the membership function and $\mu(x)$ for $x \in A$ is the membership strength of $x$ to $A$. Letting $\mathcal{V} = \{1,\ldots,N\}$ and $\mathcal{E} = \mathcal{V} \times \mathcal{V}$, one then interprets the weights as defining the membership strength of edges to the graph, i.e., one defines the "fuzzy graph" $\mathcal{G}^{(k)} = (\mathcal{V},\mathcal{E},\mu^{(k)})$ where $\mu^{(k)}(l,l') = w_{ll'}^{(k)}$ is the membership strength of edge $(l,l')$ to $\mathcal{G}^{(k)}$.

This is, of course, just a reinterpretation of weighted graphs in terms of fuzzy sets, but it allows one to combine the collection $(\mathcal{G}^{(k)}, k = 1,\ldots,N)$ using simple fuzzy sets operations, namely, defining the combined (fuzzy) graph $\mathcal{G} = (\mathcal{V},\mathcal{E},\mu)$ with

$$(\mathcal{E},\mu) = \bigcup_{k=1}^{N}(\mathcal{E},\mu^{(k)})$$

being the fuzzy union of the edge sets. There are, in fuzzy logic, multiple ways to define set unions [85], and the one selected for UMAP define $(A,\mu) \cup (A',\mu') = (A \cup A', \nu)$ with $\nu(x) = \mu(x) + \mu'(x) - \mu(x)\mu'(x)$ ($\mu(x)$ and $\mu'(x)$ being defined as 0 is $x \notin A$ or $x \notin A'$ respectively). In UMAP, each edge $\mu^{(k)}(l,l')$ is non-zero only is $k = l$ or $l'$ so that

$$\mu(l,l') = w_{ll'}^{(l)} + w_{ll'}^{(l')} - w_{ll'}^{(l)}w_{ll'}^{(l')}.$$

This defines an input fuzzy graph structure on $\{1,\ldots,N\}$ that serves as target for an optimized similar structured associated with the representation $\boldsymbol{y} = (y_1,\ldots,y_N)$. This representation, since it is designed as a homogeneous representation of the data, provides a unique fuzzy graph $\mathcal{H}(\boldsymbol{y}) = (\mathcal{V},\mathcal{E},\nu(\cdot;\boldsymbol{y}))$ and the edge membership function is defined by $\nu(l,l';\boldsymbol{y}) = \varphi_{a,b}(y_l,y_{l'})$ with

$$\varphi_{a,b}(y,y') = \frac{1}{1 + a|y - y'|^b}.$$

The parameters $a$ and $b$ are adjusted so that $\varphi_{a,b}$ provides a differentiable approximation of the function

$$\psi_{\rho_0}(y,y') = \exp(-\max(0,|y - y'| - \rho_0))$$

where $\rho_0$ is an input parameter of the algorithm. This function $\psi_{\rho_0}$ takes the same form as the membership function defined for local graphs $\mathcal{G}^{(k)}$, and its replacement by $\varphi_{a,b}$ makes possible the use of gradient-based methods for the determination of the optimal $\boldsymbol{y}$ ($\psi_{\rho_0}$ is not differentiable everywhere).

The representation $\boldsymbol{y}$ is optimized by minimizing the "fuzzy set cross-entropy"

$$C(\mu\|\nu(\cdot,\boldsymbol{y})) = \sum_{(k,l)\in\mathcal{E}}\left(\mu(k,l)\log\frac{\mu(k,l)}{\nu(k,l|\boldsymbol{y})} + (1 - \mu(k,l))\log\frac{1 - \mu(k,l)}{1 - \nu(k,l|\boldsymbol{y})}\right)$$

or, equivalently, maximizing (using, for short, $\varphi = \varphi_{a,b}$)

$$F(\mathbf{y}) = \sum_{(k,l)\in\mathcal{E}} (\mu(k,l)\log\nu(k,l|\mathbf{y}) + (1 - \mu(k,l))\log(1 - \nu(k,l|\mathbf{y})))$$

$$= \sum_{(k,l)\in\mathcal{E}} (\mu(k,l)\log\varphi(y_k,y_l) + (1 - \mu(k,l))\log(1 - \varphi(y_k,y_l)))$$

Note the important simplification compared to the similar function $F$ is t-SNE, in that the logarithm of a potentially large sum is avoided. We have

$$\partial_{y_k}F(\mathbf{y}) = 2\sum_{l=1}^{N}\mu(k,l)\partial_{y_k}\log\varphi(y_k,y_l) + 2\sum_{l=1}^{N}(1 - \mu(k,l))\partial_{y_k}\log(1 - \varphi(y_k,y_l))$$

$$= 2\sum_{l=1}^{N}\mu(k,l)\partial_{y_k}\log\frac{\varphi(y_k,y_l)}{1 - \varphi(y_k,y_l)} + 2\sum_{l=1}^{N}\partial_{y_k}\log(1 - \varphi(y_k,y_l)).$$

The optimization can be implemented using stochastic gradient ascent. Introduce random variables $\xi_{kl}$ and $\xi'_{kl}$ both taking value in $\{0,1\}$, all independent of each other and such that $P(\xi_{kl} = 1) = \mu_{kl}$ and $P(\xi'_{kl} = 1) = \epsilon$. Define

$$H_k(\mathbf{y},\xi,\xi') = 2\sum_{l=1}^{N}\xi_{kl}\partial_{y_k}\log\frac{\varphi(y_k,y_l)}{1 - \varphi(y_k,y_l)} + 2c_k\sum_{l=1}^{N}\sum_{l'=1}^{N}\xi_{kl}\xi'_{kl'}\partial_{y_k}\log(1 - \varphi(y_k,y_{l'})).$$

Then, if one takes $c_k = 1/(\epsilon\sum_l\mu(k,l))$ one has

$$E(H_k(\mathbf{y},\xi,\xi')) = \partial_{y_k}F(\mathbf{y}).$$

This corresponds to SGA iterations in which:

(1) Each edge $(k,l)$ is selected with probability $\mu(k,l)$ (which are zero for unless $k$ and $l$ are neighbors);

(2) If $(k,l)$ is selected, one selects an additional edges $(k,l')$ each with probability $\epsilon$.

Letting $l_1,\dots,l_m$ be the number of edges selected, $y_k$ is updated according to

$$y_k \leftarrow y_k + 2\gamma\left(\partial_{y_k}\log\frac{\varphi(y_k,y_l)}{1 - \varphi(y_k,y_l)} + c_k\sum_{j=1}^{m}\partial_{y_k}\log(1 - \varphi(y_k,y_{l'}))\right).$$

**Remark 21.3** If one prefers using probability rather than fuzzy set theory, the graphs $\mathcal{G}^{(k)}$ may also be interpreted as random graphs in which edges are added independently from each other and each edge $(l,l')$ is drawn with probability $\mu^{(k)}(l,l')$. The

combined graph $\mathcal{G}$ is then the random graph in which $(l,l')$ is present if and only if it is in at least one of the $\mathcal{G}^{(k)}$ and the objective function $C$ coincides with the KL divergence between this random graph and the random graph similarly defined for $y$.

However, this fuzzy/random graph formulation of UMAP—which corresponds to current practical implementations—is only a special case of the theoretical construction made in McInnes et al. [130] which builds on the theory of (fuzzy) simplicial sets and their representation of metric spaces. We refer the interested reader to this reference, which requires a mathematical background beyond the scope of these notes. ♦

# Chapter 22

# Generalization Bounds

We provide, in this chapter, an introduction to some theoretical aspects of statistical (or machine) learning, mostly focusing on the derivation of "generalization bounds" that provide high-probability guarantees on the generalization error of predictors using training data. While these bounds are not always of practical use, because making them small in realistic situations would require an enormous amount of training data, their derivations and the form they take for specific model classes bring important insight on the structure of the learning problem, and help understand why some methods may perform well while others do not.

## 22.1   Notation

We here recall some notation introduced in chapter 5. We consider a pair of random variables $(X, Y)$, with $X : \Omega \to \mathcal{R}_X$ and $Y : \Omega \to \mathcal{G}$. Regression problems correspond to $\mathcal{R}_Y = \mathbb{R}$ (or $\mathbb{R}^q$ if multivariate) and classification to $\mathcal{R}_Y$ being a finite set. A predictor is a function $f : \mathcal{R}_X \to \mathcal{R}_Y$. The general prediction problem is to find such a predictor within a class of functions, denoted $\mathcal{F}$, minimizing the prediction (or generalization error)

$$R(f) = \mathbb{E}(r(Y, f(X)))$$

where $r : \mathcal{R}_Y \times \mathcal{R}_Y \to [0. + \infty)$ is a risk function.

A training set is a family $T = ((x_1, y_1), \ldots, (x_N, y_N)) \in (\mathcal{R}_X \times \mathcal{R}_Y)^N$, the set $\mathcal{T}$ of all possible training sets therefore being the set of all finite sequences in $\mathcal{R}_X \times \mathcal{R}_Y$. A training algorithm can then be seen as a function $\mathcal{A} : \mathcal{T} \to \mathcal{F}$ which associates to each training set $T$ a function $\mathcal{A}(T) = \hat{f}_T$.

Given $T \in \mathcal{T}$, The training set error associated to a function $f \in \mathcal{F}$ is

$$\hat{R}_T(f) = \frac{1}{|T|} \sum_{(x,y) \in T} r(y, f(x)))$$

and the in-sample error associated to a learning algorithm is the function $T \mapsto \mathcal{E}_T \stackrel{\Delta}{=} \hat{R}_T(\hat{f}_T)$. Fixing the size ($N$) of $T$, one also considers the random variable $\mathbb{T}$ with values in $\mathcal{T}$ distributed as an $N$-sample of the distribution of $(X, Y)$.

A good learning algorithm should be such that the generalization error $R(\hat{f}_T)$ is small, at least in average (i.e., $E(R(\hat{f}_\mathbb{T}))$ is small). Our main goal in this chapter is to describe generalization bounds trying to find upper-bounds for $R(\hat{f}_T)$ based on $\mathcal{E}_T$ and properties of the function class $\mathcal{F}$. These bounds will reflect the bias-variance trade-off, in that, even though large function classes provide smaller in-sample errors, they will also induce a large additive term in the upper-bound, accounting for the "variance" associated to the class.

**Remark 22.1** Both variables $X$ and $Y$ are assumed to be random in the previous setting, but there are often situations when one of them is "more random" than the other. Randomness in $Y$ is associated to measurement errors, or ambiguity in the decision. Randomness in $X$ more generally relates to the issue of sampling a dataset in a large dimensional space. In some cases, $Y$ is not random at all: for example, in object recognition, the question of assigning categories for images such as those depicted in fig. 22.1 has a quasi-deterministic answer. Sometimes, it is $X$ who is not random, for example when observing noisy signals where $X$ is a deterministic discretization of a time interval and $Y$ is some function of $X$ perturbed by noise. ♦

## 22.2   Penalty-based Methods and Minimum Description Length

### 22.2.1   Akaike's information criterion

We make a computation under the following assumptions. We assume a regression model $Y = f_\theta(X) + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and $f$ is some function parametrized by $\theta \in \mathbb{R}^m$. We also assume that the true distribution is actually covered by this model and represented by a parameter $\theta_0$. Let $\hat{\theta}_T$ denote the parameter estimated by least squares using a training set $T$, and denote for short $\hat{f}_T = f_{\hat{\theta}_T}$.

The in-sample error is

$$\mathcal{E}_T = \frac{1}{N} \sum_{k=1}^{N} (y_k - \hat{f}_T(x_k))^2.$$

Figure 22.1: Images extracted from the PASCAL challenge 2007 dataset [70], in which categories must be associated with images. There is little ambiguity on correct answers based on observing the image, i.e., little randomness in the variable $Y$.

We want to compare the training-set-averaged prediction error and the average in-sample error, namely compute the error bias

$$\Delta_N = \mathbb{E}(R(f_\mathbb{T})) - \mathbb{E}(\mathcal{E}_\mathbb{T}).$$

Write

$$\Delta_N = \mathbb{E}(R(f_\mathbb{T})) - R(f_{\theta_0}) + R(f_{\theta_0}) - E(\mathcal{E}_\mathbb{T}).$$

We make a heuristic argument to evaluate $\Delta_N$. We can use the fact that $\hat{\theta}_T$ minimizes the empirical error and write

$$\frac{1}{N} \sum_{k=1}^{N} (Y_k - f_{\theta_0}(X_k))^2 = \mathcal{E}_\mathbb{T} + \sigma^2 (\hat{\theta}_\mathbb{T} - \theta_0)^T J_\mathbb{T} (\hat{\theta}_\mathbb{T} - \theta_0) + o(|\hat{\theta}_\mathbb{T} - \theta_0|^2)$$

with

$$J_T = \frac{1}{2\sigma^2 N} \sum_{k=1}^{N} \partial_\theta^2 ((y_k - f_\theta(x_k))^2)_{|\theta - \hat{\theta}_T},$$

which is an $m$ by $m$ symmetric matrix.

Now, using the fact that $\theta_0$ minimizes the mean square error (since $f_{\theta_0}(x) = \mathbb{E}(Y|X = x)$), we can write, for any $T$:

$$R(f_T) = R(f_{\theta_0}) + \sigma^2 (\hat{\theta}_T - \theta_0)^T I (\hat{\theta}_T - \theta_0) + o(|\hat{\theta}_T - \theta_0|^2)$$

with

$$I = \frac{1}{2\sigma^2} \mathbb{E}(\partial_\theta^2 (Y - f_\theta(X))^2_{|\theta - \theta_0}).$$

As a consequence, we can write (taking expectations in both Taylor expansions)

$$\Delta_N = \sigma^2 \mathbb{E}\left((\hat{\theta}_\mathbb{T} - \theta_0)^T J_\mathbb{T}(\hat{\theta}_\mathbb{T} - \theta_0)\right) + \sigma^2 \mathbb{E}\left((\hat{\theta}_\mathbb{T} - \theta_0)^T I(\hat{\theta}_\mathbb{T} - \theta_0)\right) + o(\mathbb{E}(|\hat{\theta}_\mathbb{T} - \theta_0|^2)).$$

(We skip hypotheses and justification for the analysis of the residual term.)

We now note that, because we are assuming a Gaussian noise, and that the true data distribution belongs to the parametrized family, the least-square estimator is also a maximum likelihood estimator. Indeed, the likelihood of the data is

$$\frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{k=1}^N (Y_k - f_\theta(X_k))^2\right) \prod_{k=1}^N \varphi_X(X_k)$$

where $\varphi_X$ is the p.d.f. of $X$ and does not depend on the unknown parameter.

We can therefore apply classical results from mathematical statistics [194]. Under some mild smoothness assumptions on the mapping $\theta \mapsto f_\theta$, $\hat{\theta}_\mathbb{T}$ converges to $\theta_0$ in probability when $N$ tends to infinity, the matrix $J_\mathbb{T}$ converges to $I$, which is the model's Fisher information matrix, and $\sqrt{N}(\hat{\theta}_\mathbb{T} - \theta_0)$ converges in distribution to a Gaussian $\mathcal{N}(0, I^{-1})$. This implies that both $N(\hat{\theta}_\mathbb{T} - \theta_0)^T J_\mathbb{T}(\hat{\theta}_\mathbb{T} - \theta_0)$ and $N(\hat{\theta}_\mathbb{T} - \theta_0)^T I(\hat{\theta}_\mathbb{T} - \theta_0)$ converge to a chi-square distribution with $m$ degrees of freedom, whose expectation is $m$, which indicates that $\Delta_N$ has order $2\sigma^2 m/N$.

This analysis can be used to develop model selection rules, in which one chooses between models of dimensions $k_1 < k_2 < \cdots < k_q = m$ (e.g., by truncating the last coordinates of $X$). The rule suggested by the previous computation is to select $j$ minimizing

$$\mathcal{E}_T^{(j)}(\hat{f}_T) + \frac{2\sigma^2 k_j}{N},$$

where $\mathcal{E}^{(j)}$ is the in-sample error computed using the $k_j$-dimensional model. This is an example of a penalty-based method, using the so-called Akaike's information criterion (AIC) [2].

### 22.2.2 Bayesian information criterion and minimum description length

Other penalty-based methods are more size-averse and replace the constant, 2, in AIC by a function of $N = |T|$, for example $\log N$. Such a change can be justified by a Bayesian analysis, yielding the Bayesian information criterion (BIC) [174]. The approach in this case is not based on an evaluation of the error, but on an asymptotic

estimation of the posterior distribution resulting from a Bayesian model selection principle. Like in the previous section, we content ourselves with a heuristic discussion.

Let us consider a statistical model parametrized by $\theta \in \Theta$, where $\Theta$ is an open convex subset of $\mathbb{R}^m$ with p.d.f. given by

$$f(z;\theta) = \exp(\theta^T U(z) - C(\theta)),$$

with $U : \mathbb{R}^d \to \mathbb{R}^m$ and $z = (x,y)$. We are given a family of sub-models represented by $\mathcal{M}_1,\ldots,\mathcal{M}_q$, where, for each $j$, $\mathcal{M}_j$ is the intersection of $\Theta$ with a $k_j$-dimensional affine subspace of $\mathbb{R}^m$. We are also given a prior distribution for $\theta$ in which a sub-model is first chosen, with probabilities $\alpha_1,\ldots,\alpha_q$, and given that, say, $\mathcal{M}_j$ is selected, $\theta \in \mathcal{M}_j$ is chosen with a probability distribution with density $\varphi_j$ with respect to Lebesgue's measure on $\mathcal{M}_j$ (denoted $dm_j$). Given training data $T = (z_1,\ldots,z_N)$, Bayesian model selection consists in choosing the model $\mathcal{M}_j$ where $j$ maximizes the posterior log-likelihood

$$\mu(\mathcal{M}_j|T) = \log \int_{\mathbb{R}^m} \alpha_j e^{N(\theta^T \bar{U}_T - C(\theta))} \varphi_j dm_j(\theta)$$

where $\bar{U}_T = (U(z_1) + \cdots + U(z_N))/N$.

Consider the maximum likelihood estimator $\hat{\theta}_j$ within $\mathcal{M}_j$, maximizing $\ell(\theta, \bar{U}_T) = \theta^T \bar{U}_T - C(\theta)$ over $\mathcal{M}_j$. Then one has

$$\ell(\theta, \bar{U}_T) = \ell(\hat{\theta}_j, \bar{U}_T) + \frac{1}{2}(\theta - \hat{\theta}_j)^T \partial_\theta^2 \ell(\hat{\theta}_j, \bar{U}_T)(\theta - \hat{\theta}_j) + R_j(\theta, \hat{\theta}_j)|\theta - \hat{\theta}_j|^3$$

Note that the first derivative of $\ell$ is $\partial_\theta \ell = \bar{U} - E_\theta(U)$ where $E_\theta$ is the expectation for $f(\cdot, \theta)$. The second derivative is $-\text{var}_\theta(U)$ (showing that $\ell$ is concave) and the third derivative involves third-order moments of $U$ for $E_\theta$ and (like the second derivative) does not depend on $\bar{U}_T$. In particular, we can assume that, for any $M > 0$, there exists a constant $C_M$ such that whenever $\max(|\theta|, |\hat{\theta}_j|) \le M$, we have $R_j(\theta, \hat{\theta}_j) \le C_M$.

The law of large numbers implies that $\bar{U}_T$ converges to a limit when $N$ tends to infinity, and our assumptions imply that $\hat{\theta}_j$ converges to the parameter providing the best approximation of the distribution of $Z$ for the Kullback-Leibler divergence. In particular, with probability 1, there exists an $N$ such that $\hat{\theta}_j$ belongs to any large enough, but fixed, compact set. Moreover, the second derivative $\ell(\hat{\theta}_j)$ will also converge to a limit, $-\Sigma_j$.

For any $\epsilon > 0$, write

$$\int_{\mathbb{R}^m} \alpha_j e^{N(\theta^T \bar{U}_T - C(\theta))} \varphi_j dm_j(\theta)$$

$$= \int_{|\theta - \hat{\theta}_j| \leq \epsilon} e^{N(\theta^T \bar{U}_T - C(\theta))} \varphi_j dm_j(\theta) + \int_{|\theta - \hat{\theta}_j| \geq \epsilon} e^{N(\theta^T \bar{U}_T - C(\theta))} \varphi_j dm_j(\theta).$$

The second integral converges to 0 exponentially fast when $N$ tends to $\infty$. The first one behaves essentially like

$$\int_{\mathcal{M}_j} e^{-\frac{1}{2}N(\theta - \hat{\theta}_j)^T \Sigma_j^{-1}(\theta - \hat{\theta}_j) + \log \varphi_j(\theta)} dm_j(\theta).$$

Neglecting $\log \varphi_j(\theta)$, this integral behaves like $(2\pi \det(\Sigma_j/N))^{-1/2}$, whose logarithm is $(-k_j(\log N)/2)$ plus constant terms. As a consequence, we find that

$$\mu(\mathcal{M}_j \mid T) = \max_{\theta \in \mathcal{M}_j} \ell(\theta) - \frac{k_j}{2} \log N + \text{ bounded terms.}$$

Consider, as an example, linear regression with $Y = \beta_0 + b^T x + \sigma^2 \nu$ where $\nu$ is a standard Gaussian random variable. Assume that the distribution of $X$ is known, or, preferably, make the previous discussion conditional to $X_1, \ldots, X_N$. Let sub-models $\mathcal{M}_j$ correspond to the assumption that all but the first $k_j - 1$ coefficients of $b$ vanish. Then, up to bounded terms, the Bayesian estimator must minimize (over such parameters $b$)

$$\frac{1}{2\sigma^2} \sum_{k=1}^{N} (y_k - \beta_0 - b^T x_k)^2 + \frac{k_j}{2} \log N.$$

or

$$\mathcal{E}_T^{(j)} + \frac{k_j \sigma^2}{N} \log N.$$

We now turn to another interesting point of view, which provides the same penalty, based on maximum description length principle (MDL; Rissanen [162]) measuring the coding efficiency of a model.

Let us fix some notation. We assume that one has $q$ competing models for predicting $Y$ from $X$, for example, linear regression models based on different subsets of the explanatory variables. Denote these models $\mathcal{M}_1, \ldots, \mathcal{M}_q$. Each model will be seen, not as an assumption on the true joint distribution of $X$ and $Y$, but rather as a tool to efficiently encode the training set $((x_1, y_1), \ldots, (x_N, y_N))$. To describe MDL,

which selects the model that provides the most efficient code, we need to reintroduce a few basic concepts of information theory.

The entropy of a discrete probability $P$ over a set $\Omega$ is

$$\mathcal{H}_2(P) = -\sum_{x \in \Omega} p_x \log_2 p_x.$$

(The logarithm in base 2 is used because of the tradition of coding with bits in information theory.)

For a discrete random variable $X$, the entropy $\mathcal{H}_2(X)$ is $\mathcal{H}_2(P_X)$ where $P_X$ is the probability distribution of $X$. The relation between the entropy and coding theory is as follows: a code is a function which associates to any element $\omega \in \Omega$ a string of bits $c(\omega)$. The associated code-length is denoted $l_c(\omega)$, which is simply the number of bits in $c(\omega)$. When $P$ is a probability on $\Omega$, the efficiency of a code is measured by the average code-length:

$$E_P(l_c) = \sum_{\omega \in \Omega} l_c(\omega) P(\omega).$$

Shannon's theorem [175, 55] states that, under some conditions on the code (ensuring that any sequence of words can be recognized as soon as it is observed: one says that it is instantaneously decodable) the average code length can never be larger than the entropy of $P$. Moreover, it states that there exists codes that achieve this lower bound with no more than one bit loss, such that for all $\omega$, $l_c(\omega) \le -\log_2(P(\omega)) + 1$. These optimal codes, such as the Huffman code [55], can completely be determined from the knowledge of $P$. This allows one to interpret a probability $P$ on $\Omega$ as a tool for designing codes with code-lengths essentially equal to $(-\log_2 P)$.

This statement can be generalized to continuous random variables (replacing the discrete probability $P$ by a probability density function, say $\varphi$) if one introduces a coding precision level, denoted $\delta_0$, meaning that the decoded values may differ by no more than $\delta_0$ from the encoded ones. The result is that the optimal code-length at precision $\delta_0$ can be estimated (up to one extra bit) by $-\log_2 \varphi - \log_2 \delta_0$.

In our context, each model of the conditional distribution of $Y$ given $X$, with conditional density $\varphi(y|x)$, provides a way to encode the training set with a total code length, for $(y_1, \dots, y_N)$, of

$$-\sum_{k=1}^{N} \log_2 \varphi(y_k \mid x_k) - N \log_2 \delta_0$$

(working, as before, conditionally to $x_1, \dots, x_N$). We assume that the precision at which the data is encoded is fixed, which implies that the last term does not affect the

model choice. Now, assume a sequence of $m$ parametrized model classes, $\mathcal{M}_1,\dots,\mathcal{M}_m$ and let $\varphi(y \mid x, \theta, \mathcal{M}_j)$ denote the conditional distribution with parameter $\theta$ in the class $\mathcal{M}_j$. Within model $\mathcal{M}_j$, the optimal code length corresponds to the maximum likelihood:

$$-\sum_{k=1}^{N} \log_2 \varphi(x,y \mid \hat{\theta}_j, \mathcal{M}_j) = -\max_{\theta} \left( \sum_{k=1}^{N} \log_2 \varphi(x,y; \theta, \mathcal{M}_j) \right).$$

If the models are nested, which is often the case, the most efficient will always be the largest model, since the maximization is on a larger set. However, the minimum description length (MDL) principle uses the fact that, in order to decode the compressed data, the model, including its optimal parameters, has to be known, so that the complete code needs to include a model description. The decoding algorithm will then be: decode the model, then use it to decode the data.

So assume that a model (one of the $\mathcal{M}_j$'s) has a $k_j$-dimensional parameter $\theta$. Also assume that a probability distribution, $\pi(\theta \mid \mathcal{M}_j)$, is used to encode $\theta$. Also choose a precision level, $\delta_{ij}$, for each coordinate in $\theta$, $i = 1,\dots,k_j$. (Previously, we could consider the precision of the $y_k$, $\delta_0$, as fixed, but now, the precision level for parameters is a variable that will be optimized.) The total description length using this model now becomes

$$-\sum_{k=1}^{N} \log_2 \varphi(y_k \mid x_k; \theta, \mathcal{M}_j) - \log_2 \pi(\theta \mid \mathcal{M}_j) - \sum_{i=1}^{k_j} \log_2(\delta_{ij}).$$

Let $\hat{\theta}^{(j)}$ be the parameter that maximizes

$$L(\theta \mid \mathcal{M}_j) = \sum_{k=1}^{N} \log_2 \varphi(y_k \mid x_k; \theta, \mathcal{M}_j) + \log_2 \pi(\theta \mid \mathcal{M}_j)$$

If $\pi$ is interpreted as a prior distribution of the parameters, $\hat{\theta}^{(j)}$ is the *maximum a posteriori* Bayes estimator. We now take the correction caused by $(\delta_{ij}, i = 1,\dots,k_j)$ into account, by assuming that the $i$th coordinate in $\hat{\theta}^{(j)}$ is truncated to $-\log_2 \delta_i$ bits. Let $\overline{\theta}^{(j)}$ denote this approximation. A second-order expansion of $L(\theta|\mathcal{M}_j)$ around $\hat{\theta}^{(j)}$ yields (assuming sufficient differentiability)

$$L(\overline{\theta}^{(j)} \mid \mathcal{M}_j) = L(\hat{\theta}^{(j)} \mid \mathcal{M}_j) + \frac{1}{2}(\overline{\theta}^{(j)} - \hat{\theta}^{(j)})^T S_{\hat{\theta}^{(j)}}(\overline{\theta}^{(j)} - \hat{\theta}^{(j)}) + o(|\overline{\theta}^{(j)} - \hat{\theta}^{(j)}|^2)$$

where $S_\theta$ is the matrix of second derivatives of $L(\cdot \mid \mathcal{M}_j)$ at $\theta$. Approximating $\overline{\theta}^{(j)} - \hat{\theta}^{(j)}$ by $\delta^{(j)}$ (the $k_j$-dimensional vector with coordinates $\delta_{ij}$, $i = 1,\dots,k_j$), we see that the

precision should maximize

$$\frac{1}{2}(\delta^{(j)})^T S_{\hat{\theta}^{(j)}} \delta^{(j)} + \sum_{i=1}^{k_j} \log_2 \delta_{ij}.$$

Note that $S_{\hat{\theta}^{(j)}}$ must be negative semi-definite, since $\hat{\theta}$ is a local maximum. Assuming it is non-singular, the previous expression can be maximized and yields

$$S_{\hat{\theta}^{(j)}} \delta^{(j)} = -\frac{1}{\log 2} \frac{1}{\delta^{(j)}} \tag{22.1}$$

where $1/\delta^{(j)}$ is the vector with coordinates $(1/\delta_{ij})$.

Let us now make an asymptotic evaluation. Because $L(\theta \mid \mathcal{M}_j)$ includes a sum over $N$ independent terms, it is reasonable to assume that $S_{\hat{\theta}^{(j)}}$ has order $N$, and more precisely, that $S_{\hat{\theta}^{(j)}}/N$ has a limit. Rewrite (22.1) as

$$\frac{S_{\hat{\theta}^{(j)}}}{N} \sqrt{N} \delta^{(j)} = -\frac{1}{\log 2} \frac{1}{\sqrt{N} \delta^{(j)}}.$$

This implies that $\sqrt{N} \delta^{(j)}$ is the solution of an equation which stabilizes with $N$, and it is therefore reasonable to assume that the optimal $\delta_{ij}$ takes the form $\delta_{ij} = c_i(N \mid \mathcal{M}_j)/\sqrt{N}$, with $c_i(N \mid \mathcal{M}_j)$ converging to some limit when $N$ tends to infinity. The total cost can therefore be estimated by

$$-L(\hat{\theta}^{(j)} \mid \mathcal{M}_j) + \frac{k_j}{2} \log_2 N - \frac{k_j}{2} - \sum_{i=1}^{k_j} \log_2 c_i(N \mid \mathcal{M}_j)$$

The last two terms are $O(1)$, and can be neglected, at least when $N$ is large compared to $k_j$. The final criterion becomes the penalized likelihood

$$l_d(\theta \mid \mathcal{M}_j) = L(\theta \mid \mathcal{M}_j) - \frac{k_j}{2} \log_2 N$$

in which we see that the dimension of the model appears with a factor $\log_2 N$ as announced (one needs to normalize both terms by $N$ to compare with the previous paragraph).

## 22.3 Concentration inequalities

The discussion of the AIC was a first attempt at evaluating a prediction error. It was however done under very specific parametric assumptions, including the fact that the true distribution of the data was within the considered model class. It was, in

addition, a bias evaluation, i.e., we estimated how much, in average, the in-sample error was less than the generalization error. We would like to obtain upper bounds to the generalization error that hold with high probability, and rely as little as possible on assumptions on the true data distribution.

One of the main tools used in this context are concentration inequalities, which provide upper bounds on the various probabilities of events involving a large number of random variables. The current section provides a review of some of these inequalities.

### 22.3.1   Cramér's theorem

If $X_1, X_2, \ldots$ are independent, integrable random variables with identical distributions (to that of a random variable $X$), the law of large numbers tells us that the empirical mean $\bar{X}_N = (X_1 + \cdots + X_N)/N$ converges with probability one to $m = E(X)$. When the variables are square integrable, Chebychev's inequality provides an easy proof of the weak law of large numbers. Indeed,

$$\mathbb{P}\Big(|\bar{X}_n - m| > \epsilon\Big) \leq \frac{1}{\epsilon^2}\mathbb{E}\Big((\bar{X}_N - m)^2 \mathbf{1}_{|\bar{X}_n - m| > \epsilon}\Big) \leq \frac{\text{var}(\bar{X}_N)}{\epsilon^2} = \frac{\text{var}(X)}{N\epsilon^2}.$$

A stronger assumption on the moments of $X$ yields a stronger inequality. One says that $X$ has exponential moments if there exists $\lambda_0 > 0$ such that $\mathbb{E}(e^{\lambda_0|X|}) < \infty$. In this case, the cumulant-generating function, defined, for $\lambda \in \mathbb{R}$, by

$$M_X(\lambda) = \log \mathbb{E}(e^{\lambda X}) \in [0, +\infty], \tag{22.2}$$

is finite for $\lambda \in [-\lambda_0, \lambda_0]$.

Here are a few straightforward properties of the cumulant-generating function.

(i)  One has $M_X(0) = 0$.

(ii)  For any $a \in \mathbb{R}$, one has $M_{aX}(\lambda) = M_X(a\lambda)$.

(iii)  If $X_1$ and $X_2$ are independent variables, one also has

$$M_{X_1 + X_2}(\lambda) = M_{X_1}(\lambda) + M_{X_2}(\lambda).$$

In particular, $M_{X+a}(\lambda) = M_X(\lambda) + \lambda a$, so that $M_{X-E(X)}(\lambda) = M_X(\lambda) - \lambda E(X)$.

(iv)  Finally, Markov's inequality (which states that, for any non-negative variable $Y$, $P(Y > t) \leq E(Y)/t$) applied to $Y = e^{\lambda X}$ for $\lambda > 0$ yields

$$\mathbb{P}(X > t) = \mathbb{P}(e^{\lambda X} > e^{\lambda t}) \leq e^{M_X(\lambda) - \lambda t}. \tag{22.3}$$

(Note that this inequality is trivially true for $\lambda = 0$.)

From these properties, one can easily derive a concentration inequality for the mean of independent random variables. We have $M_{\bar{X}_N}(\lambda) = N M_X(\lambda/N)$ and applying (22.3) we get, for any $\lambda \geq 0$ and $t > 0$

$$\mathbb{P}(\bar{X}_N - m > t) \leq e^{-\lambda(m+t) + M_{\bar{X}_N}(\lambda)} = e^{-N\left(\frac{\lambda(m+t)}{N} - M_X\left(\frac{\lambda}{N}\right)\right)}$$

where the right-hand side may be infinite. Because this inequality is true for any $\lambda$, we have

$$\mathbb{P}(\bar{X}_N - m > t) \leq e^{-N M_{X,+}^*(m+t)}$$

where $M_{X,+}^*(u) = \sup_{\lambda \geq 0}(\lambda u - M_X(\lambda))$, which is non-negative since the maximized quantity vanishes for $\lambda = 0$. A symmetric computation yields

$$P(\bar{X}_N - m < -t) \leq e^{-N M_{X,-}^*(m-t)}$$

where $M_{X,-}^*(t) = \sup_{\lambda \leq 0}(\lambda t - M_X(\lambda))$, which is also non-negative.

Let

$$M_X^*(t) = \sup_{\lambda \in \mathbb{R}}(\lambda t - M_X(\lambda)) \geq 0 \tag{22.4}$$

(this is the Fenchel-Legendre transform of the cumulant generating function, sometimes called the Cramér transform of $X$). One has $M_X^*(m+t) = M_{X,+}^*(m+t)$ for $t > 0$. Indeed, because $x \mapsto e^{\lambda x}$ is convex, Jensen's inequality implies that

$$\mathbb{E}(e^{\lambda X}) \geq e^{\lambda m}$$

so that $\lambda(m+t) - M_X(\lambda) \leq \lambda t < 0$ if $\lambda < 0$. Similarly, $M_X^*(m-t) = M_{X,-}^*(m-t)$ for $t > 0$. We therefore have the following result.

**Theorem 22.2** *Let $X_1, \ldots, X_N$ be independent and identically distributed random variables. Assume that these variables are integrable and let $m = \mathbb{E}(X_1)$. Then, for all $t > 0$,*

$$\mathbb{P}(\bar{X}_N - m > t) \leq e^{-N M_X^*(m+t)}$$

*and*

$$\mathbb{P}(|\bar{X}_N - m| > t) \leq 2 e^{-N \min(M_X^*(m+t), M_X^*(m-t))}$$

The last inequality derives from

$$\mathbb{P}(|\bar{X}_N - m| > t) = \mathbb{P}(\bar{X}_N - m > t) + u P(\bar{X}_N - m < -t).$$

This is our first example of concentration inequality that shows that, when

$$\min(M_X^*(m+t), M_X^*(m-t)) > 0,$$

the probability of a deviation by $t$ at least of $\bar{X}_n$ from its mean decays exponentially fast. The derivation of the inequality above was quite easy: apply Markov's inequality in a parametrized form and optimize over the parameter. It is therefore surprising that this inequality is sharp, in the sense that a similar lower bound also holds. Even though we are not going to use it in the rest of this chapter, it is worth sketching the argument leading to this lower bound, which involves an interesting step making a change of measure.

Assume (without loss of generality) that $m = 0$ and consider $\mathbb{P}(\bar{X}_n > t)$. Assume, to simplify the discussion, that the supremum of $\lambda \mapsto \epsilon\lambda - M_X(\lambda)$ is attained at some $\lambda_t$. We have

$$\partial_\lambda M_X(\lambda) = \frac{\mathbb{E}(Xe^{\lambda X})}{\mathbb{E}(e^{\lambda X})}.$$

Let $q_\lambda(x) = \frac{e^{\lambda x}}{\mathbb{E}(e^{\lambda X})}$ and $\mathbb{P}_\lambda$ (with expectation $\mathbb{E}_\lambda$) the probability distribution on $\Omega$ with density $q_\lambda(X)$ with respect to $\mathbb{P}$, so that $\partial_\lambda M_X(\lambda) = \mathbb{E}_\lambda(X)$. We have, since $\lambda_t$ is a maximizer, $\mathbb{E}_{\lambda_t}(X) = t$. Moreover, fixing $\delta > 0$,

$$\begin{aligned}
\mathbb{P}(\bar{X}_N > t) &= \mathbb{E}(\mathbf{1}_{\bar{X}_N > t}) \\
&\geq \mathbb{E}(\mathbf{1}_{|\bar{X}_n - t - \delta| < \delta}) \\
&\geq \mathbb{E}\left(\mathbf{1}_{|\bar{X}_n - t - \delta| < \delta} e^{N\lambda\bar{X}_N - Nt - 2N\delta}\right) \\
&= e^{-N(t+2\delta)} M_X(\lambda)^N \mathbb{P}_\lambda(|\bar{X}_N - t - \delta| < \delta)
\end{aligned}$$

If one takes $\lambda = \lambda_{t+\delta}$, this implies that

$$\mathbb{P}(\bar{X}_N > t) \geq e^{-NM_X^*(t+\delta)} e^{-N\delta} \mathbb{P}_{\lambda_{t+\delta}}(|\bar{X}_N - t - \delta| < \delta).$$

By the law of large numbers (applied to $\mathbb{P}_{\lambda_{t+\delta}}$), $\mathbb{P}_{\lambda_{t+\delta}}(|\bar{X}_N - t - \delta| < \delta)$ tends to 1 when $N$ tends to infinity. This implies that the logarithmic rate of convergence to 0 of $\mathbb{P}(\bar{X}_N > t)$ is larger than $N(M_X^*(t + \delta) + \delta)$, for any $\delta > 0$, to be compared with the rate $NM_X^*(t)$ for the upper bound. In Large Deviation theory, the upper and lower bounds are often simplified by considering the limit of $\log \mathbb{P}(\bar{X}_N > t)/N$, which, in this case, is $M_X^*(t)$ (and this result is called Cramér's therorem).

While Cramér's upper bound is sharp, its computation requires an exact knowledge of the distribution of $X$, which is not a common situation. The following sections optimize the upper bound in situations where only partial information on the variable is known, such as its moments or its range. As a first example, we consider concentration of the mean for sub-Gaussian variables.

### 22.3.2  Sub-Gaussian variables

If $X$ has exponential moments, then, (applying again Markov's inequality)

$$P(|X| > x) \leq Ce^{-\lambda x}$$

for some positive constants $C$ and $\lambda$. Reducing if needed the value of $\lambda$, one can assume that $C$ takes some predetermined (larger than 1) value, say, $C = 2$, the simple argument being left to the reader. A random variable such that, for some $\lambda > 0$

$$\mathbb{P}(|X| > x) \le 2e^{-\lambda x}$$

is called sub-exponential (and this property is equivalent to $X$ having exponential moments). Similarly, one says that $X$ is sub-Gaussian if, some $\sigma > 0$,

$$\mathbb{P}(|X| > x) \le 2e^{-\frac{x^2}{2\sigma^2}}. \tag{22.5}$$

Sub-Gaussian random variables are such that $M(\lambda) < \infty$ for all $\lambda \in \mathbb{R}$. Indeed, for $\lambda > 0$

$$\begin{aligned}
\mathbb{E}(e^{\lambda|X|}) &= \int_0^\infty \mathbb{P}(e^{\lambda|X|} > z)dz \\
&= 1 + \int_1^\infty \mathbb{P}(|X| > \lambda^{-1}\log z)dz \\
&\le 1 + 2\int_1^\infty e^{-\frac{(\log z)^2}{2\sigma^2\lambda^2}}dz \\
&\le 1 + 2\int_1^\infty e^{x - \frac{x^2}{2\lambda^2\sigma^2}}dx \\
&\le 1 + 2\sqrt{2\pi}\lambda\sigma e^{\frac{\lambda^2\sigma^2}{2}}.
\end{aligned}$$

**Proposition 22.3** *Assume that $X$ is sub-Gaussian, so that (22.5) holds for some $\sigma^2 > 0$. Then, for any $t > 0$, we have*

$$\mathbb{P}(\bar{X}_n - E(X) > t) \le \left(1 + \frac{4t^2}{\sigma^2}\right)^N e^{-\frac{Nt^2}{2\sigma^2}}.$$

PROOF Let us assume, without loss of generality, that $\mathbb{E}(X) = 0$. For $\lambda > 0$, we then have
$$\mathbb{E}(e^{\lambda X}) = 1 + E(e^{\lambda X} - \lambda X - 1).$$

Let $\varphi(t) = e^t - t - 1$. We have $\varphi(t) \ge 0$ for all $t$, $\varphi(0) = 0$ and, for $z > 0$, the equation $z = \varphi(t)$ has two solutions, one positive and one negative that we will denote $g_+(z) > 0 > g_-(z)$. We have

$$\begin{aligned}
\mathbb{E}(\varphi(\lambda X)) &= \int_0^\infty \mathbb{P}(\varphi(\lambda X) > z)dz \\
&= \int_0^\infty \mathbb{P}(\lambda X > g_+(z))dz + \int_0^\infty \mathbb{P}(\lambda X < g_-(z))dz
\end{aligned}$$

The change of variable $u = g_+(z)$ in the first integral is equivalent to $u > 0$, $\varphi(u) = z$ with $dz = (e^u - 1)du$. Similarly, $u = -g_-(z)$ in the second integral gives $u > 0$, $\varphi(-u) = z$ and $dz = (1 - e^{-u})du$ so that

$$\mathbb{E}(\varphi(\lambda X)) = \int_0^\infty \mathbb{P}(\lambda X > u)(e^u - 1)du + \int_0^\infty \mathbb{P}(\lambda X < -u)(1 - e^{-u})du$$

$$\leq \int_0^\infty \mathbb{P}(\lambda|X| > u)(e^u - e^{-u})du.$$

(Using the fact that $\max(\mathbb{P}(\lambda X > u), \mathbb{P}(\lambda X < -u)) \leq \mathbb{P}(\lambda|X| > u)$.) We have

$$\int_0^\infty \mathbb{P}(\lambda|X| > u)(e^u - e^{-u})du \leq 2\int_0^{+\infty}(e^u - e^{-u})e^{-\frac{u^2}{2\lambda^2\sigma^2}}du$$

$$= 2\lambda\sigma\int_0^{+\infty}(e^{\lambda\sigma v} - e^{-\lambda\sigma v})e^{-\frac{v^2}{2}}dv$$

$$= 2\lambda\sigma e^{\frac{\lambda^2\sigma^2}{2}}\sqrt{2\pi}(\Phi(-\sigma\lambda) - \Phi(\sigma\lambda))$$

$$\leq 4\lambda^2\sigma^2 e^{\frac{\lambda^2\sigma^2}{2}}$$

where $\Phi$ is the cumulative distribution function of the standard Gaussian and we have used $\Phi(-t) - \Phi(t) \leq 2t/\sqrt{2\pi}$. We therefore have

$$M_X(\lambda) \leq \log\left(1 + 4\lambda^2\sigma^2 e^{\frac{\lambda^2\sigma^2}{2}}\right) \leq \frac{\lambda^2\sigma^2}{2} + \log(1 + 4\lambda^2\sigma^2).$$

This implies

$$M_X^*(t) = \sup_{\lambda > 0}(\lambda t - M_X(\lambda)) \geq \frac{t^2}{\sigma^2} - M_X(t/\sigma^2) \geq \frac{t^2}{2\sigma^2} - \log(1 + \frac{4t^2}{\sigma^2})$$

so that

$$\mathbb{P}(\bar{X}_n > t) \leq \left(1 + \frac{4t^2}{\sigma^2}\right)^N e^{-\frac{Nt^2}{2\sigma^2}}. \qquad \blacksquare$$

The following result allows one to control the expectation of a non-negative sub-Gaussian random variable.

**Proposition 22.4** *Let X be a non-negative random variable such that*

$$\mathbb{P}(X > t) \leq Ce^{-t^2/2\sigma^2}$$

*for some constants C and $\sigma^2$. Then,*

$$\mathbb{E}(X) \leq 3\sigma\sqrt{\log C}.$$

Proof For any $\alpha \in (1, C]$, one has

$$\min(1, Ce^{-t^2/2\sigma^2}) \leq \alpha e^{-\frac{t^2 \log \alpha}{2\sigma^2 \log C}},$$

which implies that

$$\mathbb{E}(X) = \int_0^{+\infty} \mathbb{P}(X > t)dt \leq \frac{\alpha}{2\log \alpha}\sqrt{2\pi}\sigma\sqrt{\log C}$$

Taking $\alpha = \sqrt{e}$ gives

$$\mathbb{E}(X) \leq \sqrt{\pi e}\sigma\sqrt{\log C} \leq 3\sigma\sqrt{\log C}. \qquad \blacksquare$$

### 22.3.3 Bennett's inequality

The following proposition (see [24]) provides an upper bound for $M_X(\lambda)$ as a function of $\mathbb{E}(X)$ and $\mathrm{var}(X)$ under the additional assumption that $X$ is bounded from above.

**Proposition 22.5** *Let $m = \mathbb{E}(X)$ and assume that for some constant b, one has $X \leq b$ with probability one. Then, for any $\sigma^2 > 0$ such that $\mathrm{var}(X) \leq \sigma^2$, one has*

$$\mathbb{E}(e^{\lambda X}) \leq e^{\lambda m}\left(\frac{(b-m)^2}{(b-m)^2 + \sigma^2}e^{-\frac{\lambda\sigma^2}{(b-m)}} + \frac{\sigma^2}{(b-m)^2 + \sigma^2}e^{\lambda(b-m)}\right) \qquad (22.6)$$

*for any $\lambda \geq 0$.*

Proof There is no loss of generality in assuming that $m = 0$ and $\lambda = 1$, in which case one must show that

$$\mathbb{E}(e^X) \leq \frac{b^2}{b^2 + \sigma^2}e^{-\frac{\sigma^2}{b}} + \frac{\sigma^2}{b^2 + \sigma^2}e^b \qquad (22.7)$$

if $X < b$ and $E(X^2) \leq \sigma^2$. Indeed, if this inequality is true for $m = 0$ and $\lambda = 1$, (22.6) in the general case will result from letting $X = Y/\lambda + m$ and applying the special case to $Y$.

The right-hand side of (22.7) is exactly $\mathbb{E}(e^X)$ when $X$ follows the discrete distribution $P_0$ supported by two points $x_0$ and $b$, and such that $E(X) = 0$ and $\mathbb{E}(X^2) = \sigma^2$, which requires $x_0 = -\sigma^2/b$ and $P(X = x_0) = b^2/(\sigma^2 + b^2)$.

Now consider the quadratic function $v(x) = \alpha x^2 + \beta x + \gamma$ which intersects $x \mapsto e^x$ at $x = x_0$ and $x = b$, and is tangent to it at $x = x_0$, i.e., $v(b) = e^b$ and $v(x_0) = v'(x_0) = e^{x_0}$ (this uniquely defines $v$). Then $e^x \leq v(x)$ for $x < b$, yielding

$$\mathbb{E}(e^X) \leq \alpha\sigma^2 + \gamma.$$

However, since $v(X) = e^X$ almost surely when $X \sim P_0$, this upper bound is attained and equal to that provided in (22.7). $\qquad \blacksquare$

If $F(\lambda)$ denotes the right-hand side of (22.6), we have, for $m \leq u < b$,

$$M_X^*(t) \geq \sup_{\lambda \geq 0}(\lambda u - \log F(\lambda))$$

and we now estimate this lower bound. Maximizing $\lambda y - \log F(\lambda)$ is equivalent to minimizing

$$\lambda \mapsto \frac{(b-m)^2 e^{-\frac{\lambda(\sigma^2+(u-m))}{b-m}} + \sigma^2 e^{\lambda(b-u)}}{(b-m)^2 + \sigma^2}.$$

Introduce the notation $\rho = \sigma^2/(b-m)^2$, $\mu = \lambda(b-m)$ and $x = (u-m)/(b-m)$, so that the function to minimize is

$$\mu \mapsto \frac{e^{-\mu(\rho+x)} + \rho e^{\mu(1-x)}}{1+\rho}.$$

Computing the derivative in $\mu$ and equating it to 0 gives

$$\mu = \frac{1}{1+\rho} \log \frac{\rho+x}{\rho(1-x)},$$

which is non-negative since $\rho + x - \rho(1-x) = (1+\rho)x$. For this value of $\mu$, we have

$$\frac{e^{-\mu(\rho+x)} + \rho e^{\mu(1-x)}}{\rho+1} = e^{-\mu(\rho+x)} \frac{1 + \rho e^{\mu(1+\rho)}}{\rho+1}$$

$$= e^{-\mu(\rho+x)} \frac{1 + \rho \frac{\rho+x}{\rho(1-x)}}{\rho+1}$$

$$= \frac{e^{-\mu(\rho+x)}}{1-x}$$

and

$$-\log \frac{e^{-\mu(\rho+x)} + \rho e^{\mu(1-x)}}{\rho+1} = \mu(\rho+x) + \log(1-x)$$

$$= \frac{\rho+x}{1+\rho} \log \frac{\rho+x}{\rho(1-x)} + \log(1-x)$$

$$= \frac{\rho+x}{1+\rho} \log \frac{\rho+x}{\rho} + \frac{1-x}{1+\rho} \log(1-x).$$

This provides a lower bound for $M_X^*(m + (b-m)x)$, and yields the following corollary.

**Corollary 22.6** *Assume that $X$ satisfy the conditions of proposition 22.5. Then*

$$\mathbb{P}(\bar{X}_N > m + t) \leq \exp\left(-N\left(\frac{\rho+x}{1+\rho} \log \frac{\rho+x}{\rho} + \frac{1-x}{1+\rho} \log(1-x)\right)\right) \qquad (22.8)$$

*with $x = t/(b-m)$ and $\rho = \sigma^2/(b-m)^2$.*

Bennett's inequality is sometimes stated in a slightly weaker, but simpler form [127]. Returning to the proof of proposition 22.5 and using the fact that $\log u \leq u - 1$, equation (22.7) implies

$$
\begin{aligned}
\log \mathbb{E}(e^X) &\leq \frac{b^2}{b^2 + \sigma^2} e^{-\frac{\sigma^2}{b}} + \frac{\sigma^2}{b^2 + \sigma^2} e^b - 1 \\
&= \frac{b^2}{b^2 + \sigma^2}(e^{-\frac{\sigma^2}{b}} + \frac{\sigma^2}{b} - 1) + \frac{\sigma^2}{b^2 + \sigma^2}(e^b - b - 1).
\end{aligned}
$$

We will use the following lemma.

**Lemma 22.7** *The function $\varphi : u \mapsto (e^u - u - 1)/u^2$ is non-decreasing.*

PROOF We have $\varphi'(u) = \psi(u)/u^3$ where $\psi(u) = ue^u - 2e^u + u + 2$, yielding $\psi'(u) = ue^u - e^u + 1$, $\psi''(u) = ue^u$. Therefore, $\psi'$ is has its minimum at $u = 0$ with $\psi'(0) = 0$ so that $\psi$ is increasing. Since $\psi(0) = 0$, we have $\psi(u)/u^3 \geq 0$. ∎

We therefore have

$$
\begin{aligned}
\log \mathbb{E}(e^X) &\leq \frac{b^2}{b^2 + \sigma^2}(e^{-\frac{\sigma^2}{b}} + \frac{\sigma^2}{b} - 1) + \frac{\sigma^2}{b^2 + \sigma^2}(e^b - b - 1) \\
&= \frac{b^2}{b^2 + \sigma^2} \frac{\sigma^4}{b^2} \varphi(-\sigma^2/b) + \frac{\sigma^2}{b^2 + \sigma^2} b^2 \varphi(b) \\
&\leq \left(\frac{\sigma^4}{b^2 + \sigma^2} + \frac{\sigma^2 b^2}{b^2 + \sigma^2}\right)\varphi(b) \\
&= \frac{\sigma^2}{b^2}(e^b - b - 1)
\end{aligned}
$$

This shows that

$$
\log \mathbb{E}(e^{\lambda X}) \leq \frac{\sigma^2}{b^2}(e^{\lambda b} - \lambda b - 1)
$$

and

$$
M_X^*(t) \geq \frac{\sigma^2}{b^2} \max_\lambda (\lambda b^2 t/\sigma^2 - e^{\lambda b} + \lambda b + 1) = \frac{\sigma^2}{b^2} h(bt/\sigma^2)
$$

where $h(u) = (1 + u)\log(1 + u) - u$.

We summarize this in the following corollary.

**Corollary 22.8** *Assume that $X$ satisfy the conditions of proposition 22.5. Then, for $t > 0$,*

$$
\mathbb{P}(\bar{X}_N > m + t) \leq \exp\left(-\frac{N\sigma^2}{(b - m)^2} h\left(\frac{(b - m)t}{\sigma^2}\right)\right) \tag{22.9}
$$

*where $h(u) = (1 + u)\log(1 + u) - u$.*

This estimate can be further simplified as follows. Let $g$ be such that $g''(u) = (1 + u/3)^{-3}$ and $g(0) = g'(0) = 0$, which gives $g(u) = u^2/(2 + 2u/3)$. Noting that $h''(u) = (1 + u)^{-1}$ and that $(1 + u)^{-1} \geq (1 + u/3)^{-3}$, for $u \geq 0$ we find, integrating twice, that $h(u) \geq g(u)$ for $u \geq 0$. This shows that the following upper-bound is also true:

$$\mathbb{P}(\bar{X}_N > m + t) \leq \exp\left(-\frac{Nt^2}{2\sigma^2 + 2t(b-m)/3}\right). \tag{22.10}$$

This upper bound is known as Bernstein's inequality.

**Remark 22.9** It should be clear that, in the previous discussion, one may relax the assumption that $X_1, \ldots, X_N$ are identically distributed as long as there is a common function $M$ such that $M_{X_k}(\lambda) \leq m_k + M(\lambda)$ for all $k$, with $m_k = \mathbb{E}(X_k)$. We have in this case

$$\mathbb{P}(\bar{X}_N > \bar{m}_N + t) \leq \exp(-NM^*(t))$$

with $\bar{m}_N = (m_1 + \cdots + m_N)/N$ and $M^*(t) = \sup_\lambda(\lambda t - M(\lambda))$. This remark can be, in particular, applied to the situation in which $X_1, \ldots, X_N$ satisfy the conditions of proposition 22.5 with the same constants $b$ and $\sigma^2$, yielding the same upper bound as in equation (22.8). ◆

### 22.3.4   Hoeffding's inequality

We now consider the case in which the random variables $X_1, \ldots, X_N$ are bounded from above and from below, and start with the following consequence of proposition 22.5.

**Proposition 22.10** *Let $X$ be a random variable taking values in the interval $[a, b]$. Let $m = \mathbb{E}(X)$. Then*

$$\mathbb{E}(e^{\lambda X}) \leq \frac{b-m}{b-a}e^{\lambda a} + \frac{m-a}{b-a}e^{\lambda b} \leq e^{\lambda m}e^{\frac{\lambda^2(b-a)^2}{8}} \tag{22.11}$$

*for all $\lambda \in \mathbb{R}$.*

PROOF We first note that, if $X$ takes values in $[a, b]$, then $\text{var}(X) \leq (b-m)(m-a)$ (using $\sigma^2 = (b-m)(m-a)$ in (22.6)). To prove the upper bound on the variance, introduce the function $g(x) = (x-a)(x-b)$ so that $g(x) \leq 0$ on $[a, b]$. Noting that one can write $g(x) = (x-m)^2 + (2m-a-b)(x-m) + (a-m)(b-m)$, we have

$$\mathbb{E}(g(X)) = \text{var}(X) - (b-m)(m-a) \leq 0,$$

which proves the inequality.

This shows that, if $\lambda \geq 0$, we can apply proposition 22.5 with $\sigma^2 = (b-m)(m-a)$, which provides the first inequality in (22.11). To handle the case $\lambda \leq 0$, it suffices to apply this inequality with $\tilde{\lambda} = -\lambda$, $\tilde{X} = -X$, $\tilde{a} = -b$, $\tilde{b} = -a$ and $\tilde{m} = -m$.

The second inequality, namely

$$\left(\frac{b-m}{b-a}e^{\lambda a} + \frac{m-a}{b-a}e^{\lambda b}\right) \le e^{\lambda m}e^{\frac{\lambda^2(b-a)^2}{8}}$$

requires a little additional work. Letting $u = (m-a)/(b-a)$, $\alpha = \lambda(b-a)$ and taking logarithms, we need to prove that

$$\log(1 - u + ue^\alpha) - u\alpha \le \frac{\alpha^2}{8}$$

Let $f(\alpha)$ denote the difference between the right-hand side and left-hand side. Then $f(0) = 0$,

$$f'(\alpha) = \frac{\alpha}{4} - \frac{ue^\alpha}{1 - u + ue^\alpha} + u,$$

(so that $f'(0) = 0$) and

$$f''(\alpha) = \frac{1}{4} - \frac{u(1-u)e^\alpha}{(1 - u + ue^\alpha)^2}.$$

For positive numbers $x = 1 - u$ and $y = ue^\alpha$, one has $(x + y)^2 \ge 4xy$, which shows that $f''(\alpha) \ge 0$. This proves that $f'$ is non-decreasing with $f'(0) = 0$, proving that $f$ is minimized at $\alpha = 0$, so that $f(\alpha) \ge 0$ as needed. ∎

We can then deduce the following theorem [92].

**Corollary 22.11 (Hoeffding Inequality)** *If $X_1, \ldots, X_N$ are independent, taking values, respectively, in intervals of length, $c_1, \ldots, c_N$ and $Y = X_1 + \cdots + X_N$, then*

$$\mathbb{P}(Y > \mathbb{E}(Y) + t) \le \exp\left(-\frac{2t^2}{|c|^2}\right) \tag{22.12}$$

*and*

$$\mathbb{P}(Y < \mathbb{E}(Y) - t) \le \exp\left(-\frac{2t^2}{|c|^2}\right) \tag{22.13}$$

*where $|c|^2 = \sum_{k=1}^N c_k^2$.*

PROOF We have, by proposition 22.10, for any $\lambda > 0$

$$\mathbb{P}(Y > \mathbb{E}(Y) + t) \le e^{-\left(\lambda t - \sum_{k=1}^N M_{X_k}(\lambda)\right)} \le e^{-(\lambda t - \frac{\lambda^2}{8}|c|^2)}$$

The upper bound is minimized for $\lambda = 4t/|c|^2$, yielding (22.12). Equation (22.13) is obtained by applying (22.12) to $-X$. ∎

An important special case of this inequality is when $X_1, \ldots, X_N$ are i.i.d. taking values in an interval of length $\delta$. Then

$$\mathbb{P}(\bar{X}_N > \mathbb{E}(X) + t) \leq \exp\left(-\frac{2Nt^2}{\delta^2}\right). \tag{22.14}$$

This inequality is obtained after applying Hoeffding's inequality to $X_1/N, \ldots, X_N/N$, therefore taking $c_1 = \cdots = c_N = \delta/N$ and $|c|^2 = \delta^2/N$.

### 22.3.5   McDiarmid's inequality

One can relax the assumption that the random variables $X_1, \ldots, X_N$ are independent and only assume that these variables behave like "martingale increments," as stated in the following proposition [59].

**Proposition 22.12** *Let $X_1, \ldots, X_N$, $Z_1, \ldots, Z_N$ be two sequences of $N$ random variables such that*

$$\mathbb{E}(Z_k \mid X_1, Z_1, \ldots, X_{k-1}, Z_{k-1}) = m_k$$

*is constant and $|Z_k - m_k| \leq c_k$ for some constants $c_1, \ldots, c_N$. Then*

$$\mathbb{P}(Y > \mathbb{E}(Y) + t) \leq e^{-2t^2/|c|^2}$$

*with $Y = Z_1 + \cdots + Z_N$ and $|c|^2 = \sum_{k=1}^{N} c_k^2$.*

PROOF Proposition 22.10 applied to the conditional distribution implies that, for $\lambda \geq 0$:

$$\log \mathbb{E}(e^{\lambda(Z_k - m_k)} \mid X_1, Z_1, \ldots, X_{k-1}, Z_{k-1}) \leq \log \mathbb{E}(e^{\lambda|Z_k - m_k|} \mid X_1, Z_1 \ldots, X_{k-1}, Z_{k-1}) \leq \frac{\lambda^2 c_k^2}{8}.$$

Let $S_k = \sum_{j=1}^{k} (Z_j - m_j)$. Then

$$\mathbb{E}(e^{\lambda S_k}) = \mathbb{E}(e^{\lambda S_{k-1}} E(e^{\lambda(Z_k - m_k)} \mid X_1, Z_1, \ldots, X_{k-1}, Z_{k-1})) \leq e^{\frac{\lambda^2 c_k^2}{8}} \mathbb{E}(e^{\lambda S_{k-1}})$$

so that

$$\mathbb{E}(e^{\lambda S_N}) \leq e^{\frac{\lambda^2}{8} \sum_{k=1}^{N} c_k^2}$$

and the result follows from Markov's inequality optimized over $\lambda$.          ∎

We will use this proposition to prove the "bounded difference," or McDiarmid's inequality.

**Theorem 22.13 (McDiarmid's inequality)** *Let $X_1, \ldots, X_N$ be independent random variables and $g : \mathbb{R}^N \to \mathbb{R}$ a function such that there exists $c_1, \ldots, c_N$ such that*

$$|g(x_1, \ldots, x_{k-1}, x_k, x_{k+1}, \ldots, x_N) - g(x_1, \ldots, x_{k-1}, \tilde{x}_k, x_{k+1}, \ldots, x_N)| \leq c_k \qquad (22.15)$$

*for all $k = 1, \ldots, N$ and $x_1, \ldots, x_{k-1}, x_k, \tilde{x}_k, x_{k+1}, \ldots, x_N$. Then*

$$\mathbb{P}\left(g(X_1, \ldots, X_N) > \mathbb{E}(g(X_1, \ldots, X_N)) + t\right) \leq e^{-2t^2/|c|^2}$$

*with $|c|^2 = c_1^2 + \cdots + c_N^2$.*

PROOF Let $m = \mathbb{E}(g(X_1, \ldots, X_N))$. Let $Z_0 = 0$,

$$Y_k = \mathbb{E}(g(X_1, \ldots, X_N) \mid X_1, \ldots, X_k) - m$$

and $Z_k = Y_k - Y_{k-1}$. Note that $Z_k$ is a function of $X_1, \ldots, X_k$ and can therefore be omitted from the conditional expectation given $(X_1, Z_1, \ldots, X_{k-1}, Z_{k-1})$.

We have $\mathbb{E}(Y_k) = 0$ and $\mathbb{E}(Y_k \mid X_1, \ldots, X_{k-1}) = Y_{k-1}$ so that $\mathbb{E}(Z_k \mid X_1, \ldots, X_{k-1}) = 0$. Because the variables are independent, we have, letting $\tilde{X}_1, \ldots, \tilde{X}_N$ be independent copies of $X_1, \ldots, X_N$,

$$Z_k = \mathbb{E}(g(X_1, \ldots, X_{k-1}, X_k, \tilde{X}_{k+1}, \ldots, \tilde{X}_N) \mid X_1, \ldots, X_k)$$
$$- \mathbb{E}(g(X_1, \ldots, X_{k-2}, \tilde{X}_{k-1}, \tilde{X}_k, \ldots, \tilde{X}_N) \mid X_1, \ldots, X_{k-1}).$$

For fixed $X_1, \ldots, X_{k-1}$, (22.15) implies that $Z_k$ varies in an interval of length $c_k$ at most (whose bounds depend on $X_1, \ldots, X_{k-1}$) so that $|Z_k - E(Z_k)| \leq c_k$. Proposition 22.12 implies that

$$\mathbb{P}(Z_1 + \cdots + Z_N \geq t) \leq e^{-2t^2/|c|^2},$$

which concludes the proof since

$$Z_1 + \cdots + Z_N = g(X_1, \ldots, X_N) - \mathbb{E}(g(X_1, \ldots, X_N)). \qquad \blacksquare$$

### 22.3.6 Boucheron-Lugosi-Massart inequality

The following result [38], that we state without proof, extends on the same idea.

**Theorem 22.14** *Let $X_1, \ldots, X_N$ be independent random variables. Let*

$$Z = g(X_1, \ldots, X_N)$$

*with $g : \mathbb{R}^N \to [0, +\infty)$ and for $k = 1, \ldots, N$,*

$$Z_k = g_k(X_1, \ldots, X_{k-1}, X_{k+1}, \ldots, X_N)$$

*with $g_k : \mathbb{R}^{N-1} \to \mathbb{R}$. Assume that, for all $k = 1, \ldots, N$, one has $0 \leq Z - Z_k \leq 1$ and that*

$$\sum_{k=1}^{N}(Z - Z_k) \leq Z.$$

*Then*

$$\mathbb{P}(Z - \mathbb{E}(Z) > t) \leq \exp(-\mathbb{E}(Z)h(t/\mathbb{E}(Z))) \leq \exp\left(-\frac{t^2}{2\mathbb{E}(Z) + 2t/3}\right)$$

*where $h(u) = (1 + u)\log(1 + u) - u$. Moreover, for $t < \mathbb{E}(Z)$,*

$$\mathbb{P}(Z - \mathbb{E}(Z) < -t) \leq \exp(-\mathbb{E}(Z)h(-t/\mathbb{E}(Z))) \leq \exp(-t^2/2\mathbb{E}(Z)).$$

*Finally, for all $\lambda \in \mathbb{R}$*

$$\log \mathbb{E}(e^{\lambda(Z - \mathbb{E}(Z))}) \leq \mathbb{E}(Z)(e^\lambda - \lambda - 1). \tag{22.16}$$

## 22.4   Bounding the empirical error with the VC-dimension

### 22.4.1   Introduction

Section 22.3 provides some of the most important inequalities used to evaluate the deviation of various combinations of independent random variables (e.g., their empirical mean) from their expectations (the reader may refer to Ledoux and Talagrand [117], Devroye et al. [60], Talagrand [188], Dembo and Zeitouni [59], Vershynin [199] and other textbooks on the subject for further developments).

We now return to the problem of estimating the generalization error based on training data. For a given predictor $f$, concentration bounds allow us to control the probability

$$\mathbb{P}(R(f) - \hat{R}_{\mathbb{T}}(f) > t)$$

where

$$R(f) = \mathbb{E}(r(Y, f(X))$$

and

$$\hat{R}_T(f) = \frac{1}{N}\sum_{k=1}^{N} r(y_k, f(x_k))$$

for a training set $T = (x_1, y_1, \ldots, x_N, y_N)$.

If this probability is small, then $R(f) \leq \hat{R}_{\mathbb{T}}(f) + t$ with high probability, providing a likely upper bound to the generalization error of $f$. For example, if $r$ is the 0–1

loss in a classification problem, Hoeffding's inequality implies, for training sets of size $N$,

$$\mathbb{P}(R(f) - \hat{R}_{\mathbb{T}}(f) > t) \le e^{-2Nt^2}.$$

Now corollary 22.11 does not hold if we replace $f$ by $\hat{f}_{\mathbb{T}}$, i.e., if $f$ is estimated from the training set $\mathbb{T}$, which is, unfortunately, the situation we are interested in. Before addressing this problem, we point out that this inequality does apply to the case in which $f = \hat{f}_{\mathbb{T}_0}$ where $\mathbb{T}_0$ is another training set, independent from $\mathbb{T}$, so that

$$\mathbb{P}(R(\hat{f}_{\mathbb{T}_0}) - \hat{R}_{\mathbb{T}}(\hat{f}_{\mathbb{T}_0}) > t) \le e^{-2Nt^2},$$

which is proved by writing

$$\mathbb{P}(R(\hat{f}_{\mathbb{T}_0}) - \hat{R}_{\mathbb{T}}(\hat{f}_{\mathbb{T}_0}) > t) = \mathbb{E}(\mathbb{P}(R(\hat{f}_T) - \hat{R}_{\mathbb{T}}(\hat{f}_T) \mid > t \mid \mathbb{T}_0 = T)).$$

In this situation, the empirical risk is computed on a test or validation set ($\mathbb{T}$) independent of the set used to estimate $f$ ($\mathbb{T}_0$).

If one does not have a test set, and $\hat{f}_{\mathbb{T}}$ is optimized over a set $\mathcal{F}$ of possible predictors, one can rarely do much better than starting from a variation of the trivial upper bound

$$\mathbb{P}(R(\hat{f}_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}} > t) \le \mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t\right)$$

(with $\mathcal{E}_{\mathbb{T}} = \hat{R}_{\mathbb{T}}(\hat{f}_{\mathbb{T}})$) and the concentration inequalities discussed in section 22.3 need to be extended to provide upper bounds to the right-hand side.

**Remark 22.15** Computing supremums of functions over non countable sets may bring some issues regarding measurability. To avoid complications, we will always assume, when computing supremums over infinite sets, that such supremums can be reduced to maximizations over finite sets, i.e., when considering $\sup_{f \in \mathcal{F}} \Phi(f)$ for some function $\Phi$, we will assume that there exists a nested sequence of finite subsets $\mathcal{F}_n \subset \mathcal{F}$ such that

$$\sup\{\Phi(f) : f \in \mathcal{F}\} = \lim_{n \to \infty} \sup\{\Phi(f) : f \in \mathcal{F}_n\}. \tag{22.17}$$

This is true, for example, when $\mathcal{F}$ has a topology that admits a countable dense subset, with respect to which $\Phi$ is continuous. ♦

When $\mathcal{F}$ is a finite set, one can use a "union bound" with

$$\mathbb{P}(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t) \le \sum_{f \in \mathcal{F}} \mathbb{P}(R(f) - \mathcal{E}_{\mathbb{T}}(f) > t) \le |\mathcal{F}| \max_{f \in \mathcal{F}} \mathbb{P}(R(f) - \mathcal{E}_{\mathbb{T}}(f) > t).$$

Such bounds cannot be applied to the typical case in which $\mathcal{F}$ is infinite, and is likely to provide very poor estimates even when $\mathcal{F}$ is finite, but $|\mathcal{F}|$ is large. However, all proofs of concentration inequalities applied to such supremums require using a union bound at some point, often after considerable preparatory work. Union bounds will in particular appear in conjunction with the Vapnik-Chervonenkis dimension that we now discuss.

### 22.4.2   Vapnik's theorem

We consider a classification problem with two classes, 0 and 1, and therefore let $\mathcal{F}$ be a set of binary functions, i.e., taking values in $\{0, 1\}$. We also assume that the risk function $r$ takes values in the interval $[0, 1]$ (using, for example, the 0–1 loss). Let

$$U(t) = \mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t\right). \tag{22.18}$$

A fundamental theorem of Vapnik provides an estimate of $U(t)$ based on the number of possible ways to split a training set of $2N$ points into two classes using functions in $\mathcal{F}$. The rest of this section is devoted to a discussion of this result and related notions.

If $A$ is a finite subset of $\mathcal{R}$, we let $\mathcal{F}(A)$ denote the set $\{f_{|A} : f \in \mathcal{F}\}$ of restrictions of elements of $\mathcal{F}$ to the set $A$. As a convention, we let $\mathcal{F}(\emptyset) = \{f_\emptyset\}$, containing the so-called empty function. Since $\mathcal{F}$ only contains binary functions, we have $|\mathcal{F}(A)| \leq 2^{|A|}$. If $x_1, \ldots, x_M \in \mathcal{R}$, we let, with a slight abuse of notation,

$$\mathcal{F}(x_1, \ldots, x_M) = \mathcal{F}(A)$$

where $A = \{x_i, i = 1, \ldots, M\}$. This provides the number of possible splits of a training set $T = (x_1, \ldots, x_M)$ using classifiers in $\mathcal{F}$. Fixing in this section a random variable $X$, we let

$$S_{\mathcal{F}}(M) = E(|\mathcal{F}(X_1, \ldots, X_M)|)$$

where the expectation is taken over all $M$ i.i.d. realizations from $X$. We also let

$$S_{\mathcal{F}}^*(M) = \max\{|\mathcal{F}(A)| : A \subset \mathcal{R}, |A| \leq M\}.$$

The following theorem controls $U$ in (22.18) in terms of $S_{\mathcal{F}}$.

**Theorem 22.16 (Vapnik)**  *With the notation above, one has, for $t \geq \sqrt{2/N}$:*

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t\right) \leq 2S_{\mathcal{F}}(2N)e^{-Nt^2/8}, \tag{22.19}$$

*which implies that, with probability at least $1 - \delta$, we have*

$$\forall f \in \mathcal{F} : R(f) \le \mathcal{E}_{\mathbb{T}}(f)) + \sqrt{\frac{8}{N}\left(\log S_{\mathcal{F}}(N) + \log\frac{2}{\delta}\right)} \tag{22.20}$$

(The requirement that $t \ge \sqrt{2/N}$ does not really reduce the range of applicability of (22.19), since, for $t \le \sqrt{2/N}$, the upper bound in that equation is typically much larger than 1.)

PROOF We first show that the problem can be symmetrized with the inequality, valid if $Nt^2 \ge 2$,

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) \ge t\right) \le 2\mathbb{P}\left(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f)) \ge \frac{t}{2}\right) \tag{22.21}$$

in which $\mathbb{T}'$ is a second training set (independent of $\mathbb{T}$) with $N$ samples also. In view of assumption (22.17), there is no loss of generality in assuming that $\mathcal{F}$ is finite. Associate to any training set $T$, a classifier $f_T \in \mathcal{F}$ maximizing $R(f_T) - \mathcal{E}(f_T)$. One then has

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f)) \ge \frac{t}{2}\right) \ge \mathbb{P}\left((\mathcal{E}_{\mathbb{T}'}(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(f_{\mathbb{T}})) \ge \frac{t}{2}\right)$$

$$\ge \mathbb{P}\left((R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}'}(f_{\mathbb{T}}) \le \frac{t}{2} \text{ and } R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(f_T)) \ge t\right)$$

$$= \mathbb{E}\left(\mathbf{1}_{R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(f_{\mathbb{T}})) \ge t}\mathbb{P}\left(R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}'}(f_{\mathbb{T}}) \le \frac{t}{2} \mid \mathbb{T}\right)\right)$$

Conditional to $\mathbb{T}$, $\mathcal{E}_{\mathbb{T}'}(f_{\mathbb{T}})$ is the average of $M$ i.i.d. Bernoulli random variables, with variance bounded from above by $1/4$ and

$$\mathbb{P}\left(R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}'}(f_{\mathbb{T}}) \le \frac{t}{2} \mid \mathbb{T}\right) \ge 1 - \frac{1/4}{Nt^2/4} \ge \frac{1}{2}.$$

It follows that

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f)) > \frac{t}{2}\right)$$

$$\ge \frac{1}{2}\mathbb{P}\left(R(f_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(f_{\mathbb{T}})) \ge t\right) = \frac{1}{2}\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_T(f)) \ge t\right).$$

This justifies (22.21).

Now consider a family of independent Rademacher random variables $\xi_1, \ldots, \xi_N$, also independent of $\mathbb{T}$ and $\mathbb{T}'$, taking values $-1$ and $+1$ with equal probability. By symmetry,

$$\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f)) = \sup_{f \in \mathcal{F}}\sum_{k=1}^{N}(r(Y_k, f(X_k)) - r(Y_k', f(X_k')))/N$$

has the same distribution as

$$\sup_{f \in \mathcal{F}} \sum_{k=1}^{N} \xi_k (r(Y_k, f(X_k)) - r(Y_k', f(X_k')))/N .$$

Now, there are at most $|\mathcal{F}(X_1, \ldots, X_N, X_1', \ldots, X_N')|$ different sets of coefficients in front of $\xi_1, \ldots, \xi_N$ in the above sum when $f$ varies in $\mathcal{F}$, so that, conditioning on $\mathbb{T}, \mathbb{T}'$ and taking a union bound , we have

$$\mathbb{P}\left( \sup_{f \in \mathcal{F}} \sum_{k=1}^{N} \xi_k (r(Y_k, f(X_k)) - r(Y_k', f(X_k')))/N \geq t/2 \,\Big|\, \mathbb{T}, \mathbb{T}' \right)$$

$$\leq |\mathcal{F}(X_1, \ldots, X_N, X_1', \ldots, X_N')|$$

$$\sup_{f \in \mathcal{F}} \mathbb{P}\left( \sum_{k=1}^{N} \xi_k (r(Y_k, f(X_k)) - r(Y_k', f(X_k')))/N \geq t/2 \,\Big|\, \mathbb{T}, \mathbb{T}' \right)$$

The variables $\xi_k (r(Y_k, f(X_k)) - r(Y_k', f(X_k')))$ are centered and belong to the interval $[-1, 1]$, which has length 2, so that Hoeffding's inequality implies

$$\mathbb{P}\left( \sum_{k=1}^{N} \xi_k (r(Y_k, f(X_k)) - r(Y_k', f(X_k')))/N \geq t/2 \,\Big|\, \mathbb{T}, \mathbb{T}' \right) \leq e^{-2N(t/2)^2/4} = e^{-Nt^2/8}$$

and taking expectation over $\mathbb{T}$ and $\mathbb{T}'$ yields

$$\mathbb{P}\left( \sup_{f \in \mathcal{F}} (R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t \right) = 2 S_{\mathcal{F}}(2N) e^{-Nt^2/8}.$$

Equation (22.20) is then obtained from letting $\delta = 2 S_{\mathcal{F}}(2N) e^{-Nt^2/8}$ so that $t = \sqrt{\frac{8}{N} \log \frac{2 S_{\mathcal{F}}(2N)}{\delta}}$ with $R(f) \leq \mathcal{E}_{\mathbb{T}}(f) + t$ for all $f$ with probability $1 - \delta$ or more.          ∎

### 22.4.3   VC dimension

To obtain a practical bound, the quantity $S_{\mathcal{F}}(2N)$, or its upper-bound $S_{\mathcal{F}}^*(2N)$, needs to be estimated. We prove below an important property of $S_{\mathcal{F}}^*$, namely that, either $S_{\mathcal{F}}^*(M) = 2^M$ for all $M$, or there exists an $M_0$ for which $S_{\mathcal{F}}^*(M_0) < 2^{M_0}$, and taking $M_0$ to be the largest one for which an equality occurs, $S_{\mathcal{F}}^*(M)$ has order $M^{M_0}$ for all $M \geq M_0$. This motivates the following definition of the VC-dimension of the model class.

**Definition 22.17** *The Vapnik-Chervonenkis dimension (or VC dimension) of the model class $\mathcal{F}$ is*

$$VC\text{-}dim(\mathcal{F}) = \max\{M : S_{\mathcal{F}}^*(M) = 2^M\}.$$

*(where the infimum of an empty set is $+\infty$).*

**Remark 22.18** If, for a finite set $A \subset \mathcal{R}$, one has $|\mathcal{F}(A)| = 2^{|A|}$, one says that *A is shattered by* $\mathcal{F}$. So $VC\text{-}dim(\mathcal{F})$ is the largest integer $M$ such that there exists a set of cardinality $M$ in $\mathcal{R}$ that is shattered by $\mathcal{F}$. ◆

We now evaluate the growth of $S_{\mathcal{F}}^*(M)$ in terms of the VC-dimension, starting with the following lemma, which states that, if $A$ is a finite subset of $\mathcal{R}$, there are at least $|\mathcal{F}(A)|$ subsets of $A$ that are shattered by $\mathcal{F}$.

**Lemma 22.19 (Pajor)** *Let $A$ be a finite subset of $\mathcal{R}$. Then*

$$|\mathcal{F}(A)| \leq |\{B \subset A : |\mathcal{F}(B)| = 2^B\}|.$$

Proof The statement holds for $A = \emptyset$, for which $|\mathcal{F}_\emptyset| = 1 = 2^0$. For $|A| = 1$, the upper-bound is either 1 if $|\mathcal{F}(A)| = 1$, or 2 if $|\mathcal{F}(A)| = 2$, and the collection of sets $B \subset A$ such that $|\mathcal{F}(B)| = 2^B$ is $\{\emptyset\}$ in the first case and $\{\emptyset, A\}$ in the second one. So, the statement is true for $|A| = 0$ or 1.

Proceeding by induction, assume that the result is true if $|A| \leq N$, and consider a set $A'$ with $|A'| = N + 1$. Assume that $|\mathcal{F}(A')| \geq 2$ (otherwise there is nothing to prove), which implies that there exists $x \in A'$ such that $|\mathcal{F}(x)| = 2$. Take such an $x$ and write $A' = A \cup \{x\}$ with $x \notin A$. Let

$$\mathcal{F}_0 = \{f \in \mathcal{F} : f(x) = 0\} \text{ and } \mathcal{F}_1 = \{f \in \mathcal{F} : f(x) = 1\}.$$

Since $\mathcal{F}_0 \cap \mathcal{F}_1 = \emptyset$, we have

$$|\mathcal{F}(A')| = |\mathcal{F}_0(A')| + |\mathcal{F}_1(A')|.$$

Since $f(x)$ is constant on $\mathcal{F}_0$ (resp. $\mathcal{F}_1$), we have $|\mathcal{F}_0(A')| = |\mathcal{F}_0(A)|$ (resp. $|\mathcal{F}_1(A')| = |\mathcal{F}_1(A)|$), and the induction hypothesis implies

$$|\mathcal{F}(A')| \leq |\{B \subset A : |\mathcal{F}_0(B)| = 2^B\}| + |\{B \subset A : |\mathcal{F}_1(B)| = 2^B\}|$$
$$= |\{B \subset A : |\mathcal{F}_0(B)| = 2^B \text{ or } |\mathcal{F}_0(B)| = 2^B\}|$$
$$+ |\{B \subset A : |\mathcal{F}_0(B)| = |\mathcal{F}_1(B)| = 2^B\}|.$$

If $B \subset A$ is shattered by $\mathcal{F}_0$ or $\mathcal{F}_1$, it is obviously shattered by $\mathcal{F}$. Moreover, if $B$ is shattered by both, then $B \cup \{x\}$ is shattered by $\mathcal{F}$. The upper bound in the equation above is therefore less than the total number of sets shattered by $\mathcal{F}$, which proves the lemma. ∎

From this lemma, it results that if $VC\text{-}dim(\mathcal{F}) = D < \infty$, then $S_{\mathcal{F}}^*(M)$ is bounded by the total number of subsets of cardinality $D$ or less in a set of cardinality $M$. This provides the following result, which implies that the term in front of the exponential in (22.18) grows polynomially in $N$ if $\mathcal{F}$ have finite VC-dimension.

**Proposition 22.20 (Sauer-Shelah's lemma)** *If $D$ is the VC-dimension of $\mathcal{F}$, then, for $N \geq D$,*

$$S_{\mathcal{F}}^*(N) \leq \left(\frac{eN}{D}\right)^D .$$

PROOF  Pajor's lemma implies that

$$S_{\mathcal{F}}^*(N) \leq \sum_{k=0}^{D} \binom{N}{k}$$

and the statement of the proposition derives from the standard upper bound

$$\sum_{k=0}^{D} \binom{N}{k} \leq \left(\frac{eN}{D}\right)^D$$

that we now justify for completeness. We have

$$\binom{N}{k} = \frac{N!}{(N-k)!k!} \leq \frac{N^k}{k!} \leq \frac{N^D}{D^D} \frac{D^k}{k!}$$

if $k \leq D \leq N$. This yields

$$\sum_{k=0}^{D} \binom{N}{k} \leq \frac{N^D}{D^D} \sum_{k=0}^{D} \frac{D^k}{k!} \leq \frac{N^D e^D}{D^D}$$

as required.                                                                        ∎

We can therefore state a corollary to theorem 22.16 for model classes with finite VC-dimension.

**Corollary 22.21** *Assume that $VC\text{-}dim(\mathcal{F}) = D < \infty$. Then, for $t \geq \sqrt{2/N}$ and $N \geq D$,*

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t\right) \leq 2\left(\frac{2eN}{D}\right)^D e^{-Nt^2/8}. \tag{22.22}$$

*and*

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) \leq \sqrt{\frac{8}{N}} \sqrt{D \log \frac{eN}{D} + \log \frac{2}{\delta}}\right) \geq 1 - \delta. \tag{22.23}$$

### 22.4.4  Examples

The following result provides the VC-dimension of the collection of linear classifiers.

**Proposition 22.22** *Let* $\mathcal{R} = \mathbb{R}^d$ *and* $\mathcal{F} = \left\{ x \mapsto \text{sign}(a_0 + b^T x) : \beta_0 \in \mathbb{R}, b \in \mathbb{R}^d \right\}$. *Then*

$$VC\text{-}dim(\mathcal{F}) = d + 1.$$

PROOF  Let us show that no set of $d+2$ points can be shattered by $\mathcal{F}$. Use the notation $\tilde{x} = (1, x^T)^T$ and $\beta = (a_0, b^T)^T$, and consider $d + 2$ points $x_1, \ldots, x_{d+2}$. Then $\tilde{x}_1, \ldots, \tilde{x}_{d+2}$ are linearly dependent and one of them, say, $\tilde{x}_{d+2}$ can be expressed as a linear combination of the others. Write

$$\tilde{x}_{d+2} = \sum_{k=1}^{d+1} \alpha_k \tilde{x}_k.$$

Then there is no function $f \in \mathcal{F}$ (taking the form $\tilde{x} \mapsto \text{sign}(\beta \tilde{x})$) that maps $(x_1, \ldots, x_{d+2})$ to $(\text{sign}(\alpha_1), \ldots, \text{sign}(\alpha_{d+1}), -1)$ (where the definition of $\text{sign}(0) = \pm 1$ is indifferent), since any such function satisfies

$$\beta^T \tilde{x}_{d+2} = \sum_{k=1}^{d+1} \alpha_k \beta^T \tilde{x}_k > 0.$$

This proves $VC\text{-}dim(\mathcal{F}) < d + 2$. To prove that $VC\text{-}dim(\mathcal{F}) = d + 1$, it suffices to exhibit a set of $d + 1$ vectors in $\mathbb{R}^d$ that can be shattered by $\mathcal{F}$. Choose $x_1, \ldots, x_{d+1}$ such that $\tilde{x}_1, \ldots, \tilde{x}_{d+1}$ are linearly independent (for example $x_i = \sum_{k=1}^{i-1} e_i$, where $(e_1, \ldots, e_d)$ is the canonical basis of $\mathbb{R}^d$). This linear independence implies that, for any vector $\alpha = (\alpha_1, \ldots, \alpha_{d+1})^T \in \mathbb{R}^{d+1}$, there exists a vector $\beta \in \mathbb{R}^{d+1}$ such that $\tilde{x}_i^T \beta = \alpha_i$ for all $i = 1, \ldots, d + 1$. This shows that any combination of signs for $\tilde{x}_i^T \beta$ can be achieved, so that $(x_1, \ldots, x_{d+1})$ is shattered.  ∎

Upper-bounds on VC dimensions of more complex models have also been proposed in the literature. As an example, the following theorem, that we provide without proof, considers feed-forward neural networks with piecewise linear units (such as ReLU, see chapter 11). This theorem is a special case of Theorem 7 in Bartlett et al. [21], in which the more general case of networks with piecewise polynomial units is provided. Given integers $L, U_1, \ldots, U_L$ and $W_1, \ldots, W_L$, define the function class

$$\mathcal{F}(L, (U_i), (W_i), p)$$

that consists of feed-forward neural networks with $L$ layers, $U_i$ piecewise linear computational units with less than $p$ pieces in the $i$th layer, and such that the total number of parameters involved in layers $1, 2, \ldots, j$ is less than $W_j$.

**Theorem 22.23**

$$VC\text{-}dim(\mathcal{F}(L,(U_i),(W_i),p)) = O(\bar{L}W_L \log(pU)).$$

*where $U = U_1 + \cdots + U_L$ and*

$$\bar{L} = \frac{1}{W_L} \sum_{j=1}^{L} W_j.$$

Note that $p = 2$ for ReLU networks. Theorem 7 in Bartlett et al. [21] also provides a more explicit upper bound, namely

$$VC\text{-}dim(\mathcal{F}(L,(U_i),(W),p)) \leq L + \bar{L}W_L \log_2\left(4ep \sum_{i=1}^{L} iU_i \log_2\left(\sum_{i=1}^{L}(2epiU_i)\right)\right).$$

### 22.4.5   Data-based estimates

Approximations of the shattering numbers can be computed using training data. One can, in particular, prove a concentration inequality [38] on $\log S_{\mathcal{F}}(X_1,\ldots,X_N)$, which may in turn be used to estimate $\log(S_{\mathcal{F}}(2N))$. In the following, we let $\mathcal{H}_{VC}(N,\mathcal{F})$ denote the expectation of $\log S_{\mathcal{F}}(X_1,\ldots,X_N)$. It is often referred to as the VC entropy of $\mathcal{F}$.

**Theorem 22.24** *One has, letting $\mathcal{H}_{VC} = \mathcal{H}_{VC}(N,\mathcal{F})$:*

$$\mathbb{P}(\log S_{\mathcal{F}}(X_1,\ldots,X_N) \geq \mathcal{H}_{VC} + t) \leq \exp\left(-\frac{t^2}{2\mathcal{H}_{VC} + 2t/3}\right)$$

*and*

$$\mathbb{P}(\log S_{\mathcal{F}}(X_1,\ldots,X_N) \leq \mathcal{H}_{VC} - t) \leq \exp\left(-\frac{t^2}{2\mathcal{H}_{VC}}\right)$$

PROOF We show that the random variable $Z = \log_2 S_{\mathcal{F}}(X_1,\ldots,X_N)$ satisfies the assumptions of theorem 22.14, with

$$Z_k = \log_2 S_{\mathcal{F}}(X_1,\ldots,X_{k-1},X_{k+1},\ldots,X_N).$$

Clearly, $0 \leq Z$, $0 \leq Z - Z_k \leq 1$, because one can do no more than double $S_{\mathcal{F}}$ by adding one point. We need to show that

$$\sum_{k=1}^{N}(Z - Z_k) \leq Z. \tag{22.24}$$

Note that $Z$ is the base-two entropy of the uniform distribution, $\pi$, on the set

$$\mathcal{F}(X_1,\ldots,X_N) \subset \{-1,1\}^N.$$

We will use the following lemma.

**Lemma 22.25** *Let A be a finite set and $\psi$ a probability distribution on $A^N$. Let $\psi_k$ be its marginal when the kth variable is removed. Then:*

$$\sum_{k=1}^{N} \mathcal{H}_2(\psi_k) - (N-1)\mathcal{H}_2(\psi) \geq 0. \tag{22.25}$$

This lemma is a special case of a collection of results on non-negative entropy measures developed in Han [86], and we provide a direct proof below for completeness.

Given the lemma, let $\pi_k$ denote the marginal distribution of $\pi$ when the $k$th variable is removed, i.e.,

$$\pi_k(\epsilon_1, \ldots \epsilon_{k-1}, \epsilon_{k+1}, \ldots, \epsilon_N)$$
$$= \pi(\epsilon_1, \ldots \epsilon_{k-1}, -1, \epsilon_{k+1}, \ldots, \epsilon_N) + \pi(\epsilon_1, \ldots \epsilon_{k-1}, 1, \epsilon_{k+1}, \ldots, \epsilon_N).$$

We have:

$$\sum_{k=1}^{N} (\mathcal{H}_2(\pi) - \mathcal{H}_2(\pi_k)) \leq H(\pi)$$

from which (22.24) derives since $Z = \mathcal{H}_2(\pi)$ and $Z_k \geq \mathcal{H}_2(\pi_k)$. The result then follows from theorem 22.14.

We now prove lemma 22.25 by induction (this proof requires some basic notions of information theory). For convenience, introduce random variables $(\xi_1, \ldots, \xi_N)$ such that $\xi_k \in A$, with joint probability distribution given by $\psi$. Let $Y = (\xi_1, \ldots, \xi_N)$, $Y^{(k)}$ the $(N-1)$-tuple formed from $Y$ by removing $\xi_k$, $Y^{(k,l)}$ the $(N-2)$-tuple obtained by removing $\xi_k$ and $\xi_l$, etc. Inequality (22.25) can then be rewritten

$$\sum_{k=1}^{N} \mathcal{H}_2(Y^{(k)}) - (N-1)\mathcal{H}_2(Y) \geq 0.$$

This inequality is obviously true for $N = 1$, and it is true also for $N = 2$ since it gives in this case the well-known inequality $\mathcal{H}_2(Y_1, Y_2) \leq \mathcal{H}_2(Y_1) + \mathcal{H}_2(Y_2)$. Fix $M > 2$ and assume that the lemma is true for any $N < M$. To prove the statement for $N = M$, we will use the following inequality, which holds for any three random variables $U_1, U_2, U_3$:

$$\mathcal{H}_2(U_1, U_3) + \mathcal{H}_2(U_2, U_3) \geq \mathcal{H}_2(U_1, U_2, U_3) + \mathcal{H}_2(U_3).$$

This inequality is equivalent to the statement on conditional entropies that $\mathcal{H}_2(U_1, U_2 \mid U_3) \leq \mathcal{H}_2(U_1 \mid U_3) + \mathcal{H}_2(U_2 \mid U_3)$. We apply it, for given $k \neq l$, to $U_1 = Y_l$, $U_2 = Y_k$, $U_3 = Y^{(k,l)}$, yielding

$$\mathcal{H}_2(Y^{(k)}) + \mathcal{H}_2(Y^{(l)}) \geq \mathcal{H}_2(Y) + \mathcal{H}_2(Y^{(k,l)}).$$

We now sum over all pairs $k \neq l$, yielding

$$2(N-1)\sum_{k=1}^{N}\mathcal{H}_2(Y^{(k)}) \geq N(N-1)\mathcal{H}_2(Y) + \sum_{k \neq l}\mathcal{H}_2(Y^{(k,l)}).$$

We finally use the induction hypothesis to write that, for all $k$

$$\sum_{l \neq k}\mathcal{H}_2(Y^{(k,l)}) \geq (N-2)\mathcal{H}_2(Y^{(k)})$$

and obtain

$$2(N-1)\sum_{k=1}^{N}\mathcal{H}_2(Y^{(k)}) \geq N(N-1)\mathcal{H}_2(Y) + (N-2)\sum_{k=1}^{N}\mathcal{H}_2(Y^{(k)}),$$

which provides the desired result after rearranging the terms. ∎

Note that theorem 22.16 involves $S_\mathcal{F}(2N)$, with:

$$\log_2(S_\mathcal{F}(2N)) = \log_2 \mathbb{E}(S_\mathcal{F}(X_1,\ldots,X_{2N})) \geq \mathcal{H}_{VC}(2N,\mathcal{F})$$

from Jensen's inequality. This implies that the high-probability upper bound on $\mathcal{H}_{VC}(2N,\mathcal{F})$ that results from the previous theorem is not necessarily an upper bound on $\log(S_\mathcal{F}(2N))$. It is however proved in Boucheron et al. [38] that

$$\log_2 \mathbb{E}(S_\mathcal{F}(X_1,\ldots,X_{2N})) \leq \frac{1}{\log 2}\mathcal{H}_{VC}(2N,\mathcal{F})$$

also holds (as a consequence of (22.16)). A little more work (see Boucheron et al. [38]) combining theorem 22.16 and theorem 22.24 implies the following bound, which holds with probability $1 - \delta$ at least:

$$\forall f \in \mathcal{F}: R(f) \leq \mathcal{E}(f) + \sqrt{\frac{6\log S_\mathcal{F}(X_1,\ldots,X_N)}{N}} + 4\sqrt{\frac{\log(2/\delta)}{N}}.$$

## 22.5   Covering numbers and chaining

The upper bounds using the VC dimension relied on the number of different values taken by a set of functions when evaluated on a finite set, this number being used to apply a union bound. A different point of view may be applied when one relies on some notion of continuity of the family of functions on which a uniform concentration bound is needed, with respect to a given metric. This viewpoint is furthermore applicable when the sets $\mathcal{F}(X_1,\ldots,X_N)$ are infinite. To develop these tools, we will need some new concepts measuring the size of sets in a metric space.

### 22.5.1 Covering, packing and entropy numbers

**Definition 22.26** *Let $(\mathcal{G}, \rho)$ be a metric space and let $\epsilon > 0$. The $\epsilon$-covering number of $(\mathcal{G}, \rho)$. denoted $\mathcal{N}(\mathcal{G}, \rho, \epsilon)$, is the smallest integer $n$ such that there exists a subset $G \subset \mathcal{G}$ such that $|G| = n$ and $\max_{g \in \mathcal{G}} \rho(g, G) \leq \epsilon$.*

*Let $\gamma > 0$. The $\gamma$-packing number $\mathcal{M}(\mathcal{G}, \rho, \gamma)$, is the largest number $n$ such that there exists a subset $A \subset \mathcal{G}$ with cardinality $n$ such that any two distinct elements of $A$ are at distance strictly larger than $\gamma$ (such sets are called $\gamma$-nets).*

*When $\mathcal{G}$ and $\rho$ are well understood from the context, we will write simply $\mathcal{N}(\epsilon)$ and $\mathcal{M}(\gamma)$.*

**Proposition 22.27** *One has, for any $\gamma > 0$:*

$$\mathcal{M}(\mathcal{G}, \rho, 2\gamma) \leq \mathcal{N}(\mathcal{G}, \rho, \gamma) \leq \mathcal{M}(\mathcal{G}, \rho, \gamma).$$

PROOF Let $A$ be a maximal $\gamma$-net. Then, for all $x \in \mathcal{G}$, there exists $y \in A$ such that $\rho(x, y) \leq \gamma$: otherwise $A \cup \{x\}$ would also be a $\gamma - net$. This shows that $\max(\rho(x, A), x \in \mathcal{G}) \leq \gamma$ and $\mathcal{N}(\mathcal{G}, \rho, \gamma) \leq |A|$.

Conversely, let $A$ be a $2\gamma$-net. Let $G$ be an optimal $\gamma$-covering. Associate to each $y \in A$ a point $x \in G$ at distance less than $\gamma$: at least one exists because $G$ is a covering. This defines a function $f : A \to G$, which is necessarily one-to-one, because if two points in $A$ map to the same point in $G$, the distance between these two points would be less than or equal to $2\gamma$. This shows that $\mathcal{M}(\mathcal{G}, \rho, 2\gamma) \leq \mathcal{N}(\mathcal{G}, \rho, \gamma)$. ∎

The entropy numbers of $(\mathcal{G}, \rho)$, denoted, for an integer $N$, $e(\mathcal{G}, \rho, N)$ (or just $e(N)$) represent the best accuracy that can be achieved by subsets of $\mathcal{G}$ of size $N$, namely

$$e(\mathcal{G}, \rho, N) = \min_{G \subset \mathcal{G}, |G|=N} \max\{\rho(g, G) : g \in \mathcal{G}\}. \tag{22.26}$$

We have:

$$e(\mathcal{G}, \rho, N) = \inf\{\epsilon : N(\mathcal{G}, \rho, \epsilon) \leq N\} \tag{22.27a}$$

and

$$N(\mathcal{G}, \rho, \epsilon) = \min\{N : e(\mathcal{G}, \rho, N) \leq \epsilon\}. \tag{22.27b}$$

### 22.5.2 A first union bound

Let $Z$ be a random variable $Z : \Omega \to \mathcal{Z}$. We will consider a space $\mathcal{G}$ of functions $g : \mathcal{Z} \to \mathbb{R}$, such that (to simplify the discussion) $\mathbb{E}(g(Z)) = 0$ for all $g \in \mathcal{G}$. In this section, we assume that functions in $\mathcal{G}$ are bounded and let

$$\rho_\infty(g, g') = \sup_{z \in \mathcal{Z}} |g(z) - g'(z)|.$$

Assume that $\mathcal{N}(\mathcal{G}, \rho_\infty, \epsilon) < \infty$, for all $\epsilon > 0$ (which requires the set $\mathcal{G}$ to be pre-compact for the $\rho_\infty$ metric). Take $t > 0$, $0 < \epsilon < t$ and choose a set $G \subset \mathcal{G}$ such that $|G| = \mathcal{N}(\mathcal{G}, \rho_\infty, \epsilon)$. Then, using a union bound,

$$\mathbb{P}(\sup_{g \in \mathcal{G}} g(Z) \geq t) \leq \mathbb{P}(\sup_{g \in G} g(Z) \geq t - \epsilon) \tag{22.28}$$

$$\leq \mathcal{N}(\mathcal{G}, \rho_\infty, \epsilon) \sup_{g \in \mathcal{G}} \mathbb{P}(g(Z) \geq t - \epsilon).$$

Now, if each function in $\mathcal{G}$ satisfies a concentration inequality, say,

$$\mathbb{P}(g(Z) \geq u) \leq e^{-\frac{u^2}{2\mu(g)}}$$

for some $\mu(g) > 0$, then, assuming that $\mu(\mathcal{G}) \overset{\Delta}{=} \max_{g \in \mathcal{G}} \mu(g)$ is finite, we find that, for $0 < \epsilon < t$,

$$\mathbb{P}(\sup_{g \in \mathcal{G}} g(Z) \geq t) \leq \mathcal{N}(\mathcal{G}, \rho_\infty, \epsilon) e^{-\frac{(t-\epsilon)^2}{2\mu(\mathcal{G})}}.$$

We now apply this inequality to the case of binary classification, where a binary variable $Y$ is predicted by an input variable $X$, with a model class of classifiers $\mathcal{F}$ and the 0–1 loss function. If $A$ is a finite family of elements of $\mathcal{R}$, we define, for $f, f' \in \mathcal{F}$

$$\rho_A(f, f') = \frac{1}{|A|} \sum_{x \in A} \mathbf{1}_{f(x) \neq f'(x)}.$$

Let

$$\bar{\mathcal{N}}(\mathcal{F}, \epsilon, N) = E\left(\mathcal{N}(\mathcal{F}, \rho_{\{X_1,\ldots,X_N\}}, \epsilon)\right)$$

where $X_1, \ldots, X_N$ is an i.i.d. sample of $X$. We then have the following proposition.

**Proposition 22.28** *For all $\epsilon > 0$, one has*

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} (R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t\right) \leq 2\bar{\mathcal{N}}(\mathcal{F}, \epsilon/2, N) e^{-\frac{N(t/2 - \epsilon)^2}{4}}. \tag{22.29}$$

PROOF A key step in the proof of theorem 22.16, was to show that

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} (R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t\right) \leq 2\mathbb{P}\left(\sup_{f \in \mathcal{F}} \sum_{k=1}^{N} \xi_k(r(Y'_k, f(X'_k)) - r(Y_k, f(X_k))) \geq Nt/2\right). \tag{22.30}$$

where $\xi_1, \ldots, \xi_N$ are Rademacher random variables and $\mathbb{T}, \mathbb{T}'$ are two independent training sets of size $N$. We start from this inequality and bound the conditional expectation

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} \sum_{k=1}^{N} \xi_k(r(Y'_k, f(X'_k)) - r(Y_k, f(X_k))) \geq Nt/2 \,\Big|\, \mathbb{T}, \mathbb{T}'\right) \tag{22.31}$$

and therefore consider $r(Y'_k, f(X'_k)) - r(Y_k, f(X_k))$ as constants that we will denote $c_k(f)$. Since we are using a 0–1 loss, we have $c_k(f) \in \{-1, 0, 1\}$ and, for $f, f' \in \mathcal{F}$,

$$|c_k(f) - c_k(f')| \leq \mathbf{1}_{f(X_k) \neq f'(X_k)} + \mathbf{1}_{f(X'_k) \neq f'(X'_k)}. \tag{22.32}$$

Consider the random variable $Z = (\xi_1, \ldots, \xi_N)$, and let

$$\mathcal{G} = \left\{ g_f, f \in \mathcal{F} \right\}$$

with

$$g_f(\xi_1, \ldots, \xi_N) = \frac{1}{N} \sum_{k=1}^{N} c_k(f) \xi_k.$$

We have

$$\rho_\infty(g_f, g_{f'}) = \frac{1}{N} \sum_{k=1}^{N} |c_k(f) - c_k(f')|.$$

Applying Hoeffding's inequality, we have, for $u > 0$ and using the fact that $c_k \in [-1, 1]$

$$\mathbb{P}(g_f(Z) > u \mid \mathbb{T}, \mathbb{T}') \leq e^{-\frac{2Nu^2}{4}} = e^{-\frac{Nu^2}{2}}$$

and the discussion preceding the theorem yields the fact that, for any $\epsilon > 0$:

$$\mathbb{P}(\sup_{f \in \mathcal{F}} g_f(Z) > t/2 \mid \mathbb{T}, \mathbb{T}') \leq \mathcal{N}(\mathcal{G}, \epsilon, \rho_\infty) e^{-\frac{N(t/2 - \epsilon)^2}{2}}. \tag{22.33}$$

Let $A = (X_1, \ldots, X_N, X'_1, \ldots, X'_N)$ so that

$$\rho_A(f, f') = \frac{1}{2N} \sum_{k=1}^{N} \left( \mathbf{1}_{f(X_k) \neq f'(X_k)} + \mathbf{1}_{f(X'_k) \neq f'(X'_k)} \right).$$

Using (22.32), we have $\rho_\infty(g_f, g_{f'}) \leq 2\rho_A(f, f')$, which implies

$$\mathcal{N}(\mathcal{G}, \epsilon, \rho_\infty) \leq \mathcal{N}(\mathcal{F}, \epsilon/2, \rho_A).$$

Using this in (22.33) and taking the expectation in (22.31), we get

$$\mathbb{P}\left( \sup_{f \in \mathcal{F}} (R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t \right) \leq 2\bar{\mathcal{N}}(\mathcal{F}, \epsilon/2, N) e^{-\frac{N(t/2 - \epsilon)^2}{2}} \tag{22.34}$$

which is valid for all $\epsilon > 0$. ∎

One can retrieve the bound obtained in theorem 22.16 using the obvious fact that

$$\mathcal{N}(\mathcal{F}, \epsilon, \rho_A) \leq |\mathcal{F}(A)|,$$

for any $A \subset \mathcal{R}$, so that

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t\right) \leq 2\mathcal{S}(\mathcal{F}, 2N)e^{-\frac{N(t/2-\epsilon)^2}{2}}$$

for any $\epsilon > 0$, and letting $\epsilon$ go to zero,

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq t\right) \leq 2\mathcal{S}(\mathcal{F}, 2N)e^{-\frac{Nt^2}{8}}.$$

So (22.29) provides a family of equations that depend on a parameter $\epsilon$ which, in the limit $\epsilon \rightarrow 0$, includes theorem 22.16 as a particular case. For a given $N$, optimizing (22.29) over $\epsilon$ may give a better upper bound, provided one has a good way to estimate $\bar{\mathcal{N}}(\mathcal{F}, \epsilon/2, N)$ (which is, of course, far from obvious).

### 22.5.3   Evaluating covering numbers

Covering numbers can be evaluated in some simple situations. The following proposition provides an example in finite dimensions.

**Proposition 22.29** *Assume that $\mathcal{G}$ is a parametric family of functions, so that $\mathcal{G} = \{g_\theta, \theta \in \Theta\}$ where $\Theta \subset \mathbb{R}^m$. Assume also that, for some constant $C$, $\rho_\infty(g_\theta, g_{\theta'}) \leq C|\theta - \theta'|$ for all $\theta, \theta' \in \Theta$. Let $\mathcal{G}^{(M)} = \{g_\theta : \theta \in \Theta, |\theta| \leq M\}$. Then*

$$\mathcal{N}(\mathcal{G}, \rho_\infty, \epsilon) \leq \left(1 + \frac{2CM}{\epsilon}\right)^m$$

PROOF Letting $\rho$ denote the Euclidean distance in $\mathbb{R}^m$, our hypotheses imply that $\mathcal{N}(\mathcal{G}^{(M)}, \rho_\infty, \epsilon)$ is bounded by $\mathcal{N}(B_M, \rho, \epsilon/C)$ where $B_M$ is the ball with radius $M$ in $\mathbb{R}^m$. Now, if $\theta_1, \ldots, \theta_n$ is an $\alpha$-covering of $B_M$, then $\theta_1/M, \ldots, \theta_n/M$ is an $(\alpha/M)$-covering of $B_1$, which shows (together with a symmetric argument) that $\mathcal{N}(B_M, \rho, \alpha) = \mathcal{N}(B_1, \rho, \alpha/M)$ and we get

$$\mathcal{N}(\mathcal{G}^{(M)}, \rho_\infty, \epsilon) \leq \mathcal{N}(B_1, \rho, \epsilon/MC)$$

and we only need to evaluate $\mathcal{N}(B_1, \rho, \alpha)$ for $\alpha > 0$. Using proposition 22.27, one can instead evaluate $\mathcal{M}(B_1, \rho, \alpha)$. So let $A$ be an $\alpha$-net in $B_1$. Then

$$\bigcup_{x \in A} B_\rho(x, \alpha/2) \subset B_\rho(0, 1 + \alpha/2)$$

and, since the sets in the union are disjoint,

$$\sum_{x \in A} \text{volume}(B_\rho(x, \alpha/2)) = |A|\text{volume}(B_\rho(0, \alpha/2)) \leq \text{volume}(B_\rho(0, 1 + \alpha/2)).$$

Letting $C_m$ denote the volume of the unit ball in $\mathbb{R}^m$, this shows

$$|A|C_m\left(\frac{\alpha}{2}\right)^m \leq C_m\left(1 + \frac{\alpha}{2}\right)^m$$

and

$$|A| \leq \left(1 + \frac{2}{\alpha}\right)^m,$$

which concludes the proof. ∎

One can also obtain entropy number estimates in infinite dimensions. Here, we quote a result applicable to spaces of smooth functions, referring to Van der Vaart and Wellner [195] for a proof.

**Theorem 22.30** *Let $\mathcal{Z}$ be a bounded convex subset of $\mathbb{R}^d$ with non-empty interior. For $p \geq 1$ and $f \in C^p(\mathcal{Z})$, let*

$$\|f\|_{p,\infty} = \max\left\{|D^k(f(x))| : k = 0,\ldots,p, x \in \mathcal{Z}\right\}.$$

*Let $\mathcal{G}$ be the unit ball for this norm,*

$$\mathcal{G} = \left\{f \in C^p(\mathcal{Z}) : \|f\|_{p,\infty} \leq 1\right\}.$$

*Let $\mathcal{Z}^{(1)}$ be the set of all $x \in \mathbb{R}^d$ at distance less than 1 from $\mathcal{R}$.*

*Then there exists a constant $K$ depending only on $p$ and $d$ such that*

$$\log \mathcal{N}(\epsilon, \mathcal{G}, \rho_\infty) \leq K \text{volume}(\mathcal{Z}^{(1)})\left(\frac{1}{\epsilon}\right)^{d/p}$$

### 22.5.4 Chaining

The distance $\rho_\infty$ may not always be the best one to analyze the set of functions, $\mathcal{G}$. For example, if $\mathcal{G}$ is a class of functions with values in $\{-1, 1\}$, then $\rho_\infty(g, g') = 2$ unless $g = g'$. In such contexts, it is often preferable to use distances that compute average discrepancies, such as

$$\rho_p(g, g') = E(|g(Z) - g'(Z)|^p)^{1/p}, \tag{22.35}$$

for some random variable $Z$. Such distances, by definition, do not provide uniform bounds on differences between functions (that we used to write (22.28)), but can rather be used in upper-bounds on the probabilities of deviations from zero, which have to be handled somewhat differently. We here summarize a general approach called "chaining," following for this purpose the presentation made in Talagrand [189] (see also Audibert and Bousquet [15]). From now on, we assume that $(\mathcal{G}, \rho)$ is a

(pseudo-)metric space of functions $g : \mathcal{Z} \to \mathcal{R}$ and $Z$ a random variable taking values in $\mathcal{Z}$. We will make the basic assumption that, for all $g, g' \in \mathcal{G}$ and $t > 0$,

$$\mathbb{P}(|g(Z) - g'(Z)| > t) \leq 2e^{-\frac{t^2}{2\rho(g,g')^2}}.$$

Note that this assumption includes cases in which

$$\mathbb{P}(|g(Z) - g'(Z)| > t) \leq 2e^{-\frac{t^2}{2\rho(g,g')^\alpha}}.$$

for some $\alpha \in (0, 2]$, because, if $\rho$ is a distance, then so is $\rho^{\alpha/2}$ if $\alpha \leq 2$. We will also assume that $\mathbb{E}(g(Z)) = 0$ in order to avoid centering the variables at every step.

We are interested in upper bounds for $\mathbb{P}(\sup_{g \in \mathcal{G}} g(Z) > t)$. To build a chaining argument, consider a family $(G_0, G_1, \ldots)$ of subsets of $\mathcal{G}$. Assume that $|G_k| \leq N_k$ with $N_k$ chosen, for future simplicity, so that $N_{k-1} N_k \leq N_{k+1}$. For $g \in \mathcal{G}$, let $\pi_k(g)$ denote a closest point to $g$ in $G_k$. Also assume that $G_0 = \{g_0\}$ is a singleton, so that $\pi_0(g) = g_0$ for all $g \in \mathcal{G}$. (One can generally assume without harm that $0 \in \mathcal{G}$, in which case one should choose $g_0 = 0$ in the following discussion.) For $g \in G_n$, we therefore have

$$g - g_0 = \sum_{k=1}^{n} (\pi_k(g) - \pi_{k-1}(g)).$$

Let $(t_1, t_2, \ldots)$ be a sequence of numbers that will be determined later. Let

$$S_n = \max_{g \in G_n} \sum_{k=1}^{n} t_k \rho(\pi_k(g), \pi_{k-1}(g)). \tag{22.36}$$

Then, for any $t$,

$$\mathbb{P}(\sup_{g \in G_n} g(Z) - g_0(Z) > tS_n)$$
$$\leq \mathbb{P}(\exists g \in G_n, \exists k \leq n : \pi_k(g)(Z) - \pi_{k-1}(g)(Z) > tt_k \rho(\pi_k(g), \pi_{k-1}(g)))$$
$$\leq \mathbb{P}(\exists k \leq n, \exists g \in G_k, g' \in G_{k-1} : g(Z) - g'(Z) > tt_k \rho(g, g'))$$
$$\leq \sum_{k=1}^{n} N_k N_{k-1} \sup_{g \in G_k, g' \in G_{k-1}} \mathbb{P}(g(Z) - g'(Z) > tt_k \rho(g, g'))$$
$$\leq 2 \sum_{k=1}^{n} N_{k+1} e^{-\frac{t^2 t_k^2}{2}}$$

If one takes $N_k = 2^{2^k}$, which satisfies $N_k N_{k-1} = 2^{2^k + 2^{k-1}} \leq N_{k+1}$, and $t_k = 2^{k/2}$, one finds that

$$\mathbb{P}(\sup_{g \in G_n} g(Z) - g_0(Z) > tS_n) \leq 2 \sum_{k=1}^{n} 2^{2^{k+1}} e^{-2^{k-1} t^2}.$$

The upper bound converges (as a function of $n$) as soon as $t > 2\sqrt{\log 2}$. Moreover, one has

$$2\sum_{k=1}^{n} 2^{2^{k+1}} e^{-2^{k-1}t^2} = 2e^{-\frac{t^2}{2}} \sum_{k=1}^{n} e^{-2^{k-2}(t^2 - 8\log 2)} \le 2e^{-\frac{t^2}{2}} \sum_{k=1}^{\infty} e^{-2^{k-2}}$$

when $t > \sqrt{1 + 8\log 2}$. This provides a concentration bound for $\mathbb{P}(\sup_{g \in G_n} g(Z) - g_0(Z) > tS_n)$, that we may rewrite as

$$\mathbb{P}(\sup_{g \in G_n} g(Z) - g_0(Z) > t) \le C e^{-\frac{t^2}{2S_n^2}} \tag{22.37}$$

for $t > 2S_n\sqrt{\log 2}$, $C = 2\sum_{k=1}^{\infty} e^{-2^{k-2}}$ and $S_n$ given by (22.36), with $t_k = 2^{k/2}$. Moreover, we have

$$S_n = \max_{g \in G_n} \sum_{k=1}^{n} 2^{k/2} \rho(\pi_k(g), \pi_{k-1}(g))$$

$$\le \max_{g \in G_n} \sum_{k=1}^{n} 2^{k/2} (\rho(g, G_k) + \rho(g, G_{k-1}))$$

$$\le 2 \max_{g \in G_n} \sum_{k=0}^{n} 2^{k/2} \rho(g, G_k)$$

and this simpler upper bound can be used in (22.37).

We haven't made many assumptions so far on the sequence $G_0, G_1, \ldots$, beyond bounding their cardinality, but it is natural to require that they are built in order to behave like a dense subset of $\mathcal{G}$, so that

$$\lim_{n \to \infty} \max_{g \in \mathcal{G}} \rho(x, G_n) = 0. \tag{22.38}$$

Note that this requires that the set $\mathcal{G}$ is precompact for the distance $\rho$. We will also assume that

$$\lim_{n \to \infty} \sup_{g \in G_n} g(x) = \sup_{g \in \mathcal{G}} g(x). \tag{22.39}$$

Then, we have proved the following result [188].

**Theorem 22.31** *Let $G_0, G_1, \ldots$ be a family of subsets of $\mathcal{G}$ satisfying (22.38) and (22.39) and such that $G_0 = \{g_0\}$ and $|G_n| \le 2^{2^n}$ for $n \ge 0$. Let*

$$S = 2\sup_{g \in \mathcal{G}} \sum_{n=0}^{\infty} 2^{n/2} \rho(g, G_n) \tag{22.40}$$

*Then, for $t > S\sqrt{1 + 8\log 2}$,*

$$\mathbb{P}(\sup_{g \in \mathcal{G}} g(Z) - g_0(Z) > t) \leq Ce^{-\frac{t^2}{2S^2}} \tag{22.41}$$

*with $C = 2\sum_{k=1}^{\infty} e^{-2^{k-2}}$.*

The exponential rate of convergence in the right-hand side of (22.41) is the quantity $S$, and the upper bound will be improved when building the sequence $(G_0, G_1, \ldots)$ so that $S$ is as small as possible. Such an optimization for a given family of functions is however a formidable problem. It is however interesting to see (still following [188]) that theorem 22.31 implies a classical inequality in terms of what is called the metric entropy of the metric space $(\mathcal{G}, \rho)$.

### 22.5.5   Metric entropy

If $S$ is given by (22.40), we have

$$S = 2\sup\left(\sum_{n=0}^{\infty} 2^{n/2}\rho(g, G_n) : g \in \mathcal{G}\right) \leq 2\sum_{n=0}^{\infty} 2^{n/2}\sup\{\rho(g, G_n) : g \in \mathcal{G}\}$$

Take $G_n$ achieving the minimum in the entropy number $e(\mathcal{G}, \rho, 2^{2^n})$. Then, (22.41) holds with $S$ replaced by

$$\hat{S} = 2\sum_{n=0}^{\infty} 2^{n/2}e(\mathcal{G}, \rho, 2^{2^n}).$$

Consider the function

$$h(\mathcal{G}, \rho) = \int_0^{\infty} \sqrt{\log \mathcal{N}(\mathcal{G}, \rho, \epsilon)}d\epsilon, \tag{22.42}$$

which is known as *Dudley's metric entropy* of the space $(\mathcal{G}, \rho)$. We have

$$h(\mathcal{G}, \rho) = \int_0^{e(2)} \sqrt{\log \mathcal{N}(\epsilon)}d\epsilon + \sum_{n=1}^{\infty} \int_{e(2^{2^{n-1}})}^{e(2^{2^n})} \sqrt{\log \mathcal{N}(\epsilon)}d\epsilon.$$

If $\epsilon \in [e(2^{2^{n-1}}), e(2^{2^n}))$, we have $\mathcal{N}(\epsilon) > 2^{2^n}$ so that

$$h(\mathcal{G}, \rho) \geq e(2)\sqrt{\log 3} + \sum_{n=1}^{\infty} 2^{n/2}(e(2^{2^n}) - e(2^{2^{n-1}}))$$

$$\geq \left(1 - \frac{\sqrt{2}}{2}\right)\sum_{n=1}^{\infty} 2^{n/2}e(2^{2^n}).$$

Therefore,

$$\hat{S} \leq \frac{4}{2 - \sqrt{2}} h(\mathcal{G}, \rho) \leq 7 h(\mathcal{G}, \rho)$$

and this upper bound can also be used to obtain a simpler (but weaker) form of theorem 22.31.

**Remark 22.32** The covering numbers of a class $\mathcal{G}$ of binary functions $g$ with values in $\{-1, 1\}$ can be controlled by the VC dimension of the class. Here, we consider $\rho(g, g') = \mathbb{P}(g \neq g') = \rho_1(g, g')/2$. Then, the following theorem holds.

**Theorem 22.33** *Let $\mathcal{G}$ be a class of binary functions such that $D = VC\text{-}dim(\mathcal{G}) < \infty$. Then, there is a universal constant $K$ such that, for any $\epsilon \in (0, 1)$,*

$$\mathcal{N}(\mathcal{G}, \rho, \epsilon) \leq K D (4e)^D \left(\frac{1}{\epsilon}\right)^{D-1}$$

*with $\rho(g, g') = \mathbb{P}(g \neq g')$.*

We refer to Van der Vaart and Wellner [195], Theorem 2.6.4 for a proof, which is rather long and technical.                                                                              ◆

### 22.5.6 Application

We quickly show how this discussion can be turned into results applicable to the classification problem. If $\mathcal{F}$ is a function class of binary classifiers and $r$ is the risk function, one can consider the class

$$\mathcal{G} = \{(x, y) \mapsto r(y, f(x)) : f \in \mathcal{F}\}.$$

If $r$ is the 0–1 loss, we have $VC\text{-}dim(\mathcal{G}) \leq VC\text{-}dim(\mathcal{F})$. Indeed, if one considers $N$ points in $\mathcal{R} \times \{-1, 1\}$, say $(x_1, y_1, \ldots, x_N, y_N)$, then

$$\mathcal{G}(x_1, y_1, \ldots, x_N, y_N)$$
$$= \{r(1, f(x_k)) : k = 1, \ldots, N, y_k = 1\} \cup \{r(-1, f(x_k)) : k = 1, \ldots, N, y_k = -1\}.$$

If the two sets in the right-hand side are not empty, i.e., the numbers $N_{(1)}$ and $N_{(-1)}$ of $k$'s such that $y_k = 1$ or $y_k = -1$ are not zero, then

$$|\mathcal{G}(x_1, y_1, \ldots, x_N, y_N)| \leq 2^{N_{(1)}} + 2^{N_{(-1)}},$$

which is less that $2^N$ as soon as $N > 2$. So, taking $N > 2$, for $(x_1, y_1, \ldots, x_N, y_N)$ to be shattered by $\mathcal{G}$, we need $N_{(1)} = N$ or $N_{(-1)} = N$ and in this case, the inequality:

$$|\mathcal{G}(x_1, y_1, \ldots, x_N, y_N)| \leq |\mathcal{F}(x_1, \ldots, x_N)|$$

is obvious. The same inequality will be true for some $x_1,\dots,x_N$ with $N = 2$, except in the uninteresting case where $f(x) = 1$ (or $-1$) for every $x \in \mathcal{R}$.

A similar inequality holds for entropy numbers with the $\rho_1$ distance (cf. (22.35)) because

$$\mathbb{E}(|r(Y,f(X)) - r(Y,f'(X))|) \leq \mathbb{P}(f(X) \neq f'(X))$$

whenever $r$ takes values in $[0,1]$, which implies that

$$\mathcal{N}(\mathcal{G},\rho_1,\epsilon) \leq \mathcal{N}(\mathcal{F},\rho_1,\epsilon)$$

for all $\epsilon > 0$. Note however that evaluating this upper bound may still be challenging and would rely on strong assumptions on the distribution of $X$ allowing to control $\mathbb{P}(f(X) \neq f'(X))$.

We now assume that functions in $\mathcal{F}$ define "posterior probabilities" on $\mathcal{G}$. More precisely, given $\lambda \in \mathbb{R}$ we can define the probability $\pi_\lambda$ on $\{-1,1\}$ by

$$\pi_\lambda(y) = \frac{e^{\lambda y}}{e^{-\lambda} + e^\lambda}.$$

Now, if $\mathcal{F}$ is a class of *real-valued* functions, we can define the risk function

$$r(y,f(x)) = \log \frac{1}{\pi_{f(x)}(y)}.$$

Since $|\partial_\lambda \log \pi_\lambda(y)| = |y - \tanh\lambda| \leq 2$ for $y \in \{-1,1\}$, we have

$$|r(y,f(x)) - r(y,f'(x))| \leq 2|f(x) - f'(x)|$$

so that entropy numbers in $\mathcal{G}$ can be estimated from entropy numbers in $\mathcal{F}$. As an example, let $\mathcal{F}$ be a space of affine functions $x \mapsto a_0 + b^T x$, $x \in \mathbb{R}^d$. Assume that the random variable $X$ is bounded, so that one can take $\mathcal{R}$ to be an open ball centered at $0$ with radius, say, $U$. For $M > 0$, let

$$\mathcal{F}_M = \{f : x \mapsto a_0 + b^T x : |b| \leq M, |a_0| \leq UM\}.$$

The restriction $|b| \leq M$ is equivalent to using a penalty method, such as, for example, ridge logistic regression. Moreover, if $|b| \leq M$, it is natural to assume that $|a_0| \leq UM$ because otherwise $f$ would have a constant sign on $\mathcal{R}$. In this case, we get

$$\rho_\infty(r(y,f(x)),r(y,f'(x))) \leq |a_0 - a_0'| + U|b - b'|$$

and a small modification of the proof of proposition 22.29 shows that

$$\mathcal{N}(\mathcal{F},\rho_\infty,\epsilon) \leq \left(1 + \frac{4CU}{\epsilon}\right)^{d+1}$$

## 22.6 Other complexity measures

### 22.6.1 Fat-shattering and margins

VC-dimension and metric entropy are measures that control the complexity of a model class, and can therefore be evaluated a priori without observing any data. These bounds can be improved, in general, by using information derived from the training set, and, particular the classification margin that has been obtained [18].

For this discussion, we need to return to the definition of covering numbers. If $\mathcal{F}$ is a function class, $\rho_\infty$ the supremum metric on $\mathcal{F}$, $\epsilon > 0$ and $N$ is an integer, we let

$$\mathcal{N}(\mathcal{F}, \rho_\infty, \epsilon, N) = \max\{\mathcal{N}(\mathcal{F}(A), \rho_\infty, \epsilon) : A \subset \mathcal{R}, |A| = N\}$$

that we will abbreviate in $\mathcal{N}_\infty(\epsilon, N)$ when $\mathcal{F}$ is known from the context. We will assume that functions in $\mathcal{F}$ take values values in $[-1, 1]$, and we define for $\gamma \geq 0$, $y \in \{0, 1\}$, $u \in \mathbb{R}$:

$$r_\gamma(y, u) = \begin{cases} 0 \text{ if } u < -\gamma \text{ and } y = 0 \\ 0 \text{ if } u > \gamma \text{ and } y = 1 \\ 1 \text{ otherwise} \end{cases}$$

So, $r_\gamma(y, f(x))$ is equal to 0 if $f(x)$ correctly predicts $y$ with margin $\gamma$ and to 1 otherwise. We then define the classification error with margin $\gamma$ as

$$R_\gamma(f) = E(r_\gamma(Y, f(X)))$$

and, given a training set $T$ of size $N$

$$\mathcal{E}_{\gamma, T} = \frac{1}{N} \sum_{k=1}^{N} r_\gamma(y_k, f(x_k)).$$

We then have the following theorem [10].

**Theorem 22.34** *If $t \geq \sqrt{2/N}$*

$$\mathbb{P}(\sup_{f \in \mathcal{F}}(R_0(f) - \mathcal{E}_{\gamma, \mathbb{T}}(f)) > t) \leq 2\mathcal{N}_\infty(\gamma/2, 2N)e^{-Nt^2/8}, \tag{22.43}$$

*or, equivalently, with probability larger than $1 - \delta$, one has, for all $f \in \mathcal{F}$,*

$$R_0(f) - \mathcal{E}_{\gamma, \mathbb{T}}(f)) \leq \sqrt{\frac{8}{N}\left(\log \mathcal{N}_\infty(\gamma/2, 2N) + \log \frac{2}{\delta}\right)}. \tag{22.44}$$

Proof We first note that, for $Nt^2 > 2$,

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}}(R_0(f) - \mathcal{E}_{\gamma, \mathbb{T}}(f)) > t\right) \leq 2\mathbb{P}\left(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\gamma, \mathbb{T}}(f)) > \frac{\epsilon}{2}\right),$$

which is proved exactly the same way as (22.21) in theorem 22.16, and we skip the argument.

We have

$$\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\gamma,\mathbb{T}}(f) = \frac{1}{N} \sum_{k=1}^{N} (r_0(Y_k', f(X_k')) - r_\gamma(Y_k, f(X_k)))$$

and because $(X_k, Y_k)$ and $(X_k', Y_k')$ have the same distribution, $\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\gamma,\mathbb{T}}(f))$ has the same distribution as

$$\Delta_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) =$$

$$\sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{k=1}^{N} \Big( (r_0(Y_k', f(X_k')) - r_\gamma(Y_k, f(X_k)))\xi_k + (r_0(Y_k, f(X_k)) - r_\gamma(Y_k', f(X_k')))(1 - \xi_k) \Big)$$

where $\xi_1, \ldots, \xi_N$ is a sequence of Bernoulli random variables with parameter 1/2.

We now estimate $\mathbb{P}(\Delta_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) > t/2 \mid \mathbb{T}, \mathbb{T}')$ and we therefore consider $\mathbb{T}$ and $\mathbb{T}'$ as fixed. Let $F$ be a subset of $\mathcal{F}$, with cardinality $\mathcal{N}_\infty(\gamma/2, 2N)$, such that for all $f \in \mathcal{F}$ there exists an $f' \in F$ such that $|f(x) - f'(x)| \leq \gamma/2$ for all $x \in \{X_1, \ldots, X_N, X_1', \ldots, X_N'\}$. Then we claim that

$$\Delta_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) \leq \Delta'_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N)$$

where

$$\Delta'_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) = \max_{f \in F} \frac{1}{N} \sum_{k=1}^{N} (2\xi_k - 1)\Big( r_{\frac{\gamma}{2}}(Y_k', f(X_k')) - r_{\frac{\gamma}{2}}(Y_k, f(X_k)) \Big).$$

This is because, for any $(x, y) \in \mathcal{R} \times \{0, 1\}$, and $f, f'$ such that $|f(x) - f'(x)| < \gamma/2$, we have $r_0(y, f(x)) \leq r_{\gamma/2}(y, f'(x))$ and $r_{\gamma/2}(y, f'(x)) \leq r_\gamma(y, f(x))$: if an example is misclassified by $f$ (resp. $f'$) at a given margin, it must be misclassified by $f'$ (resp. $f$) at this margin plus $\gamma/2$.

Now,

$$\mathbb{P}(\Delta'_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) > \frac{t}{2})$$

$$\leq |F| \max_{f \in F} \mathbb{P}\Big( \frac{1}{N} \sum_{k=1}^{N} (2\xi_k - 1)(r_{\frac{\gamma}{2}}(Y_k', f(X_k')) - r_{\frac{\gamma}{2}}(Y_k, f(X_k))) > \frac{t}{2} \Big)$$

to which we can apply Hoeffding's inequality, yielding

$$\mathbb{P}\Big( \Delta'_{\mathbb{T},\mathbb{T}'}(\xi_1, \ldots, \xi_N) > \frac{t}{2} \Big) \leq |F| e^{-Nt^2/8}, \qquad \blacksquare$$

which concludes the proof, since, by proposition 22.27, $|F| \leq \mathcal{N}_\infty(\gamma/2, 2N)$.

In order to evaluate the covering numbers $\mathcal{N}_\infty(\epsilon, N)$ using quantities similar to VC-dimensions, a different type of set decomposition and shattering has been proposed. Following Alon et al. [4], we introduce the following notions. Recall that a family of functions $\mathcal{F} : \mathcal{R} \to \{0, 1\}$ shatters a finite set $A \subset \mathcal{R}$ if and only if $|\mathcal{F}(A)| = 2^{|A|}$. The following definitions are adapted to functions taking values in a continuous set.

**Definition 22.35** *Let $\mathcal{F}$ be a family of functions $f : \mathcal{R} \to [-1, 1]$ and $A$ a finite subset of $\mathcal{R}$.*

*(i) One says that $\mathcal{F}$ P-shatters $A$ if there exists a function $g_A : \mathcal{R} \to \mathbb{R}$ such that, for each $B \subset A$, there exists a function $f \in \mathcal{F}$ such that $f(x) \geq g_A(x)$ if $x \in B$ and $f(x) < g_A(x)$ if $x \in A \setminus B$.*

*(ii) Let $\gamma$ be a positive number. One says that $\mathcal{F}$ $P_\gamma$-shatters $A$ if there exists a function $g_A : \mathcal{R} \to \mathbb{R}$ such that, for each $B \subset A$, there exists a function $f \in \mathcal{F}$ such that $f(x) \geq g_A(x) + \gamma$ if $x \in B$ and $f(x) \leq g_A(x) - \gamma$ if $x \in A \setminus B$.*

Note that only the restriction of $g_A$ to $A$ matters in this definition. This function acts as a threshold for binary classification. More precisely, given a function $g : A \to \mathbb{R}$, one can associate to every $f \in \mathcal{F}$ the binary function $f_g$ with $f_g(x)$ equal to 1 if $f(x) \geq g(x)$ and to 0 otherwise. Letting $\mathcal{F}_g = \{f_g : f \in \mathcal{F}\}$ we see that $\mathcal{F}$ P-shatters $A$ if there exists a function $g_A$ such that $\mathcal{F}_{g_A}$ shatters $A$. The definition of $P_\gamma$-shattering introduces a margin in the definition of $f_g$ (with $f_g(x)$ equal to 1 if $f(x) \geq g(x) + \gamma$, to 0 if $f(x) \leq g(x) - \gamma$ and is ambiguous otherwise), and $A$ is $P_\gamma$-shattered by $\mathcal{F}$ if, for some $g_A$, the corresponding $\mathcal{F}_{g_A}$ shatters $A$ without ambiguities.

**Definition 22.36** *One then defines the P-dimension of $\mathcal{F}$ by*

$$P\text{-}dim(\mathcal{F}) = \max\{|A| : A \subset \mathcal{R}, \mathcal{F} \text{ P-shatters } A\},$$

*and similarly the $P_\gamma$-dimension of $\mathcal{F}$ is*

$$P_\gamma\text{-}dim(\mathcal{F}) = \max\{|A| : A \subset \mathcal{R}, \mathcal{F} \text{ } P_\gamma\text{-shatters } A\}.$$

The $P_\gamma$-dimension of $\mathcal{F}$ will replace the VC-dimension in order to control the covering numbers. More precisely, we have the following theorem [4].

**Theorem 22.37** *Let $\gamma > 0$ and assume that $\mathcal{F}$ has $P_{\gamma/4}$-dimension $D < \infty$. Then,*

$$\mathcal{N}_\infty(\gamma, N) \leq 2 \left( \frac{16N}{\gamma^2} \right)^{D \log(4eN/(D\gamma))}.$$

Proof  The proof is quite technical and relies on a combinatorial argument in which $\mathcal{F}$ is first assumed to take integer values before addressing the continuous case.

*Step 1.* We first assume that functions in $\mathcal{F}$ take values in the finite set $\{1,\ldots,r\}$ where $r$ is an integer. For the time of this proof, we introduce yet another notion of shattering called S-shattering (for strong shattering) which is essentially the same as $P_1$-shattering, except that functions $g$ are restricted to take values in $\{1,\ldots,r\}$. Let $A$ be a finite subset of $\mathcal{R}$. Given a function $g : \mathcal{R} \to \{1,\ldots,r\}$, we say that $(\mathcal{F},g)$ S-shatters $A$ if, for any $B \subset A$, there exist $f \in \mathcal{F}$ satisfying $f(x) \geq g(x) + 1$ for $x \in B$ and $f(x) \leq g(x) - 1$ if $x \in A \setminus B$. We say that $\mathcal{F}$ S-shatters $A$ if $(\mathcal{F},g)$ S-shatters $A$ for some $g$. The S-dimension of $\mathcal{F}$ is the cardinality of the largest subset of $\mathcal{R}$ that can be S-shattered and will be denoted $S\text{-}dim(\mathcal{F})$. The first, and most difficult, part of the proof is to show that, if $S\text{-}dim(\mathcal{F}) = D$, then

$$\mathcal{M}(\mathcal{F}(A),\rho_\infty,2) \leq 2(|A|r^2)^{\lceil \log_2 y \rceil}$$

with

$$y = \sum_{k=1}^{D} \binom{|A|}{k} r^k$$

and $\lceil i \rceil$ denotes the smallest integer larger than $u \in \mathbb{R}$. Here, $\mathcal{M}$ is the packing number defined in section 22.5.1.

To prove this, we can assume that $r \geq 3$, since, for $r \leq 2$, $\mathcal{M}(\mathcal{F}(A),\rho_\infty,2) = 1$ (the diameter of $\mathcal{F}$ for the $\rho_\infty$ distance is 0 or 1). Let $\mathcal{G}(A) = \{1,\ldots,r\}^A$ be the set of all functions $f : A \to \{1,\ldots,r\}$ and let

$$\mathcal{U}_A = \{F \subset \mathcal{G}(A) : \forall f, f' \in F, \exists x \in A \text{ with } |f(x) - f'(x)| \geq 2\}.$$

For $F \in \mathcal{U}_A$, let

$$\mathcal{S}_A(F) = \{(B,g) : B \subset A, B \neq \emptyset, g : B \to \{1,\ldots,r\}, (F,g) \text{ S-shatters } B\}.$$

Let $t_A(h) = \min\{|\mathcal{S}_A(F)| : F \in \mathcal{U}_A, |F| = h\}$ (where the minimum of the empty set is $+\infty$). Since we are considering in $\mathcal{U}_A$ all possible functions from $A$ to $\{1,\ldots,r\}$, it is clear that $t_A(h)$ only depends on $|A|$, and we will also denote it by $t(h,|A|)$.

Note that, by definition, if $(B,g) \in \mathcal{S}_A(F)$, and $F \subset \mathcal{F}$, then $|B| \leq D$. So, the number of elements in $\mathcal{S}_A(F)$ for such an $F$ is less or equal than the number of possible such pairs $(B,g)$, which is strictly less than $y = \sum_{k=1}^{D} \binom{|A|}{k} r^k$. So, if $t(h,|A|) \geq y$, then there cannot be any $F \subset \mathcal{F}$ in the set $\mathcal{U}_A$ and $\mathcal{M}(\mathcal{F}(A),\rho_\infty,2) < h$. The rest of the proof consists in showing that $t(h,|A|) \geq y$.

For any $n \geq 1$, we have $t(2,n) = 1$: fix $x \in A$, and $F = \{f_1, f_2\} \in \mathcal{G}$ such that $f_1(x) = 1$, $f_2(x) = 3$ and $f_1(y) = f_2(y)$ if $y \neq x$. Then only $(\{x\},g)$ is S-shattered by $F$, with $g$ such that $g(x) = 2$.

Now, assume that, for some integer $m$, $t(2mnr^2, n) < \infty$, so that there exists $F \in \mathcal{U}_A$ such that $|F| = 2mnr^2$. Arrange the elements of $F$ into $mnr^2$ pairs $\{f_i, f_i'\}$. For each such pair, there exists $x_i \in A$ such that $|f_i(x_i) - f_i'(x_i)| > 1$. Since there are at most $n$ selected $x_i$, one of them must be appearing at least $mr^2$ times. Call it $x$ and keep (and reindex) the corresponding $mr^2$ pairs, still denoted $\{f_i, f_i'\}$. Now, there are at most $r(r-1)/2$ possible distinct values for the unordered pairs $\{f_i(x), f_i'(x)\}$, so that one of them must be appearing at least $2mr^2/r(r-1) > 2m$ times. Select these functions, reindex them and exchange the role of $f_i$ and $f_i'$ if needed to obtain $2m$ pairs $\{f_i, f_i'\}$ such that $f_i(x) = k$ and $f_i'(x) = l$ for all $i$ and fixed $k, l \in \{1, \dots, r\}$ such that $k + 1 < l$. Let $F_1 = \{f_1, \dots, f_{2m}\}$ and $F_1' = \{f_1', \dots, f_{2m}'\}$. Let $A' = A \setminus \{x\}$. Then both $F_1$ and $F_1'$ belong to $\mathcal{U}_{A'}$, which implies that both $\mathcal{S}_{A'}(F_1)$ and $\mathcal{S}_{A'}(F_1')$ have cardinality at least $t(2m, n-1)$. Moreover, both sets are included in $\mathcal{S}_A(F)$, and if $(B, g) \in \mathcal{S}_{A'}(F_1) \cap \mathcal{S}_{A'}(F_1')$, then $(B \cup \{x\}, g') \in \mathcal{S}_A(F)$, with $g'(y) = g(y)$ for $y \in B$ and $g'(x) = k + 1$. This provides $2t(2m, n-1)$ elements in $\mathcal{S}_A(F)$ and shows the key inequality (which is obviously true when the left-hand side is infinite)

$$t(2mnr^2, n) \geq 2t(2m, n-1).$$

This inequality can now be used to prove by induction that for all $0 \leq k < n$, one has $t(2(nr^2)^k, n) \geq 2^k$, since

$$t(2((n+1)r^2)^{k+1}, n+1) \geq 2t(2((n+1)r^2)^k, n) \geq 2t(2(nr^2)^k, n).$$

For $k \geq n$, one has $2(nr^2)^k > r^n$, where $r^n$ is the number of functions in $\mathcal{G}(A)$, so that $t(2(nr^2)^k, n) = +\infty$. So, $t(2(nr^2)^k, n) \geq 2^k$ is valid for all $k$ and it suffices to take $k = \lceil \log_2 y \rceil$ to obtain the desired result.

*Step 2.* The next step uses a discretization scheme to extend the previous result to functions with values in $[-1, 1]$. More precisely, given $f : \mathcal{R} \to [0, 1]$, and $\eta > 0$, let

$$f^\eta(x) = \max\{k \in \mathbb{N} : 2k\eta - 1 \leq f(x)\}$$

which takes values in $\{0, \dots, r\}$ for $r = \lfloor \eta^{-1} \rfloor$. If $\mathcal{F}$ is a class of functions with values in $[-1, 1]$, define $\mathcal{F}^\eta = \{f^\eta : f \in \mathcal{F}\}$. With this notation, the following holds.

(a) For all $\gamma \leq \eta$: $S\text{-}dim(\mathcal{F}^\eta) \leq P_\gamma\text{-}dim(\mathcal{F})$

(b) For all $\epsilon \geq 4\eta$ and $A \subset \mathcal{R}$: $\mathcal{M}(\mathcal{F}(A), \rho_\infty, \epsilon) \leq \mathcal{M}_\infty(\mathcal{F}^\eta(A), \rho_\infty, 2)$.

To prove (a), assume that $\mathcal{F}^\eta$ S-shatters $A$, so that there exists $g$ such that, for all $B \subset A$, there exists $f \in \mathcal{F}$ such that $f^\eta(x) \geq g(x) + 1$ for $x \in B$ and $f^\eta(x) \leq g(x) - 1$ for $x \in A \setminus B$. Using the fact that $2\eta f^\eta(x) - 1 \leq f(x) < 2\eta f^\eta(x) + 2\eta - 1$, we get $f(x) \geq 2\eta g(x) + 2\eta - 1$ for $x \in B$ and $f(x) \leq 2\eta g(x) - 1$ for $x \in A \setminus B$. So taking $\tilde{g}(x) = 2\eta g(x) + \eta - 1$

as threshold function (which does not depend on $B$), we see that $\mathcal{F}$ $P_\gamma$-shatters $A$ if $\gamma \leq \eta$.

For (b), we deduce from the definition of $f^\eta$ that $|f^\eta(x) - \tilde{f}^\eta(x)| > (2\eta)^{-1}|f(x) - \tilde{f}(x)| - 1$ so that, if $\epsilon = 4\eta$, $|f(x) - \tilde{f}(x)| \geq \epsilon$ implies $|f^\eta(x) - \tilde{f}^\eta(x)| > 1$, or, equivalently $|f^\eta(x) - \tilde{f}^\eta(x)| \geq 2$.

*Step 3.* We can now conclude. Taking $\gamma > 0$, we have, if $|A| = N$

$$\mathcal{N}(\mathcal{F}(A), \rho_\infty, \gamma) \leq \mathcal{M}(\mathcal{F}(A), \rho_\infty, \gamma) \leq \mathcal{M}(\mathcal{F}^{\gamma/4}(A), \rho_\infty, 2) \leq 2\left(\frac{16N}{\gamma^2}\right)^{\lceil \log y \rceil}$$

with

$$y = \sum_{k=1}^{D} \binom{N}{k}(\gamma/4)^{-k} \leq (\gamma/4)^{-D} \sum_{k=1}^{D} \binom{N}{k} \leq \left(\frac{4Ne}{D\gamma}\right)^D.$$

Since the maximum of $\mathcal{N}(\mathcal{F}(A), \rho_\infty, \gamma)$ over $A$ with cardinality $N$ is $\mathcal{N}_\infty(\gamma, N)$, the proof is complete.                                                             ∎

One can use this result to evaluate margin bounds on linear classifiers with bounded data. Let $\mathcal{R}$ be the ball with radius $\Lambda$ in $\mathbb{R}^d$ and consider the model class containing all functions $f(x) = a_0 + b^T x$ with $a_0 \in [-\Lambda, \Lambda]$ and $b \in \mathbb{R}^d$, $|b| \leq 1$. Let $A = \{x_1, \ldots, x_N\}$ be a finite subset of $\mathcal{R}$. Then, $\mathcal{F}$ $P_\gamma$-shatters $A$ if and only if there exists $g_1, \ldots, g_N \in \mathbb{R}$ such that, for any sequences $\xi = (\xi_1, \ldots, \xi_N) \in \{-1, 1\}^N$, there exists $a_0^\xi \in [-\Lambda, \Lambda]$ and $b^\xi \in \mathbb{R}^d$, $|b^\xi| \leq 1$ with $\xi_k(a_0^\xi + (b^\xi)^T x_k - g_k) \geq \gamma$ for $k = 1, \ldots, N$. Summing over $N$, we find that

$$N\gamma + \sum_{k=1}^{N} g_k \xi_k \leq a_0^\xi \sum_{k=1}^{N} \xi_k + (b^\xi)^T \sum_{k=1}^{N} \xi_k x_k.$$

This shows that, for any sequence $\xi_1, \ldots, \xi_N$,

$$N\gamma + \sum_{k=1}^{N} g_k \xi_k \leq \Lambda \left| \sum_{k=1}^{N} \xi_k \right| + \left| \sum_{k=1}^{N} \xi_k x_k \right|$$

Applying the same inequality after changing the signs of $\xi_1, \ldots, \xi_N$ yields

$$N\gamma \leq N\gamma + \left| \sum_{k=1}^{N} g_k \xi_k \right| \leq \Lambda \left| \sum_{k=1}^{N} \xi_k \right| + \left| \sum_{k=1}^{N} \xi_k x_k \right|.$$

This shows, in particular, that (letting $\xi_1, \ldots, \xi_N$ be independent Rademacher random variables)

$$\mathbb{P}\left( \Lambda \left| \sum_{k=1}^{N} \xi_k \right| + \left| \sum_{k=1}^{N} \xi_k x_k \right| \geq N\gamma \right) = 1.$$

However, using the identity $(A+B)^2 \leq 2A^2 + 2B^2$, we have

$$\mathbb{P}\left(\Lambda\left|\sum_{k=1}^{N}\xi_k\right| + \left|\sum_{k=1}^{N}\xi_k x_k\right| \geq N\gamma\right) \leq P\left(2\Lambda^2\left|\sum_{k=1}^{N}\xi_k\right|^2 + 2\left|\sum_{k=1}^{N}\xi_k x_k\right|^2 \geq N^2\gamma^2\right).$$

Since

$$\mathbb{E}\left(2\Lambda^2\left|\sum_{k=1}^{N}\xi_k\right|^2 + 2\left|\sum_{k=1}^{N}\xi_k x_k\right|^2\right) = 2N\Lambda^2 + 2\sum_{k=1}^{n}|x_k|^2 \leq 4N\Lambda^2,$$

Markov's inequality implies

$$\mathbb{P}\left(2\Lambda^2\left|\sum_{k=1}^{N}\xi_k\right|^2 + 2\left|\sum_{k=1}^{N}\xi_k x_k\right|^2 \geq N^2\gamma^2\right) \leq \frac{4\Lambda^2}{N\gamma^2}.$$

We get a contradiction unless $N \leq 4\Lambda^2/\gamma^2$, which shows that $P_\gamma\text{-}dim(\mathcal{F}) \leq 4\Lambda^2/\gamma^2$. Theorem 22.37 then implies that

$$\mathcal{N}_\infty(\gamma, N) \leq 2\left(\frac{16N}{\gamma^2}\right)^{\frac{63\Lambda^2}{\gamma^2}\log\left(\frac{16eN\gamma}{\Lambda^2}\right)}$$

and this upper bound can then be plugged into equations (22.43) or (22.44) to estimate the generalization error.

Beyond the explicit expression of the upper bound, the important point in the previous argument is that the $P_\gamma$-dimension is bounded independently from the dimension $d$ of $X$ (and therefore also applies in the infinite-dimensional case). This should be compared to what we found for the VC-dimension of separating hyperplanes, which was $d+1$ (cf. proposition 22.22).

**Remark 22.38** Note that the upper-bound obtained in theorem 22.34 depends on a parameter $(\gamma)$ and the result is true for any choice of this parameter. It is tempting at this point to optimize the bound with respect to $\gamma$, but this would be a mistake since a family of events being likely does not imply that their intersection is likely too. However, with a little work, one can ensure that an intersection of slightly weaker inequalities holds. Indeed, assume that an estimate similar to (22.43) holds, in the form

$$\mathbb{P}(R_0(\hat{f}_{\mathbb{T}}) > U_{\mathbb{T}}(\gamma) + t) \leq C(\gamma)e^{-mt^2/2}$$

or, equivalently

$$\mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > U_{\mathbb{T}}(\gamma) + \sqrt{t^2 + 2\log C(\gamma)}\right) \leq e^{-mt^2/2},$$

where $U_{\mathbb{T}}(\gamma)$ depends on data and is increasing (as a function of $\gamma$), and $C(\gamma)$ is a decreasing function of $\gamma$. Consider a decreasing sequence $(\gamma_k)$ that converges to 0 (for example $\gamma_k = L2^{-k}$). Choose also an increasing function $\epsilon(\gamma)$. Then

$$\mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > \min\{U_{\mathbb{T}}(\gamma) + \sqrt{t^2 + 2\log C(\gamma) + \epsilon^2(\gamma)} : 0 \le \gamma \le L\}\right)$$
$$\le \mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > \min\{U_{\mathbb{T}}(\gamma_k) + \sqrt{t^2 + 2\log C(\gamma_{k-1}) + \epsilon^2(\gamma_k)} : k \ge 1\}\right).$$

Moreover

$$\mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > \min\{U_{\mathbb{T}}(\gamma_k) + \sqrt{t^2 + 2\log C(\gamma_{k-1}) + \epsilon^2(\gamma_k)} : k \ge 1\}\right)$$
$$\le \sum_{k=0}^{\infty} \mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > U_{\mathbb{T}}(\gamma_k) + \sqrt{t^2 + 2\log C(\gamma_{k-1}) + \epsilon(\gamma_k)}\right)$$
$$\le \sum_{k=0}^{\infty} \frac{C(\gamma_k)}{C(\gamma_{k-1})} e^{-m\epsilon^2(\gamma_k)/2 - mt^2/2}.$$

So, it suffices to choose $\epsilon(\gamma)$ so that

$$C_0 = \sum_{k=1}^{\infty} \frac{C(\gamma_k)}{C(\gamma_{k-1})} e^{-m\epsilon^2(\gamma_k)/2} < \infty$$

to ensure that

$$\mathbb{P}\left(R_0(\hat{f}_{\mathbb{T}}) > \min\{U_{\mathbb{T}}(\gamma) + \sqrt{t^2 + 2\log C(\gamma) + \epsilon^2(\gamma)} : \gamma_0 \le \gamma \le L\}\right) \le C_0 e^{-mt^2/2}.$$

For example, if $\gamma_k = L2^{-k}$, one can take

$$\epsilon(\gamma) = \sqrt{\frac{2}{m}\left(\log \frac{C(\gamma)}{C(\gamma/2)} + \log \gamma^{-1}\right)}$$

which yields $C_0 \le L$.                                          ◆

### 22.6.2  Maximum discrepancy

Let $T$ be a training set and let $T_1$ and $T_2$ form a fixed partition of the training set in two equal parts. Assume, for simplicity, that $N$ is even and that the method for selecting the two parts is deterministic, e.g., place the first half of $T$ in $T_1$ and second one in $T_2$. Following Bartlett et al. [20], one can then define the maximum discrepancy on $T$ by

$$C_T = \sup_{f \in \mathcal{F}}(\mathcal{E}_{T_1}(f) - \mathcal{E}_{T_2}(f))$$

This discrepancy measures the extent to which estimators may differ when trained on two independent half-sized training sets. For a binary classification problem, the estimation of $C_T$ can be made with the same algorithm as the initial classifier, since $\mathcal{E}_{T_1}(f) - \mathcal{E}_{T_2}(f)$ is, up to a constant, exactly the classification error for the training set in which the class labels are flipped for the data in $T_2$.

Following [20], we now discuss concentration bounds that rely on $C_T$ and start with the following Lemma.

**Lemma 22.39** *Introduce the function*

$$\Phi(T) = \sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_T(f)) - \sup_{f \in \mathcal{F}}(\mathcal{E}_{T_1}(f) - \mathcal{E}_{T_2}(f)).$$

*Then* $\mathbb{E}(\Phi(\mathbb{T})) \leq 0$.

PROOF  Note that, if $\mathbb{T}'$ is a training set, independent of $\mathbb{T}$ with identical distribution, then, for any $f_0 \in \mathcal{F}$,

$$R(f_0) - \mathcal{E}_{\mathbb{T}}(f_0) = \mathbb{E}(\mathcal{E}_{\mathbb{T}'}(f_0) - \mathcal{E}_{\mathbb{T}}(f_0) \mid \mathbb{T}) \leq \mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f)) \mid \mathbb{T})$$

so that

$$\mathbb{E}(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_T(f))) \leq \mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}'}(f) - \mathcal{E}_{\mathbb{T}}(f))).$$

Now, for a given $f$, we have $\mathcal{E}_{\mathbb{T}}(f) = \frac{1}{2}(\mathcal{E}_{\mathbb{T}_1}(f) + \mathcal{E}_{\mathbb{T}_2}(f))$ and splitting $\mathbb{T}'$ the same way, we have $\mathcal{E}_{\mathbb{T}'}(f) = \frac{1}{2}(\mathcal{E}_{\mathbb{T}_1'}(f) + \mathcal{E}_{\mathbb{T}_2'}(f))$.

We can therefore write

$$\mathbb{E}(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_T(f))) \leq \frac{1}{2}\mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}_1'}(f) - \mathcal{E}_{\mathbb{T}_1}(f)) + (\mathcal{E}_{\mathbb{T}_2'}(f) - \mathcal{E}_{\mathbb{T}_2}(f)))$$

$$\leq \frac{1}{2}\left( \mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}_1'}(f) - \mathcal{E}_{\mathbb{T}_1}(f))) + \mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}_2'}(f) - \mathcal{E}_{\mathbb{T}_2}(f))) \right)$$

$$= \mathbb{E}(\sup_{f \in \mathcal{F}}(\mathcal{E}_{\mathbb{T}_1}(f) - \mathcal{E}_{\mathbb{T}_2}(f)))$$

where we have used the fact that both $(\mathbb{T}_1', \mathbb{T}_1)$ and $(\mathbb{T}_2', \mathbb{T}_2)$ form random training sets with identical distribution to $(\mathbb{T}_1, \mathbb{T}_2)$.

This proves that $\mathbb{E}(\Phi(\mathbb{T})) \leq 0$. ∎

Using the lemma, one can write

$$\mathbb{P}(\sup_{f \in \mathcal{F}}(R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq C_{\mathbb{T}} + \epsilon) = \mathbb{P}(\Phi(\mathbb{T}) \geq \epsilon) \leq \mathbb{P}(\Phi(\mathbb{T}) - \mathbb{E}(\Phi(\mathbb{T})) \geq \epsilon).$$

One can then use McDiarmid's inequality (theorem 22.13) after noticing that, letting $z_k = (x_k, y_k)$ for $k = 1, \ldots, N$,

$$\max_{z_1, \ldots, z_N, z_k'} \left| \Phi(z_1, \ldots, z_N) - \Phi(z_1, \ldots, z_{k-1}, z_k', z_{k+1}, \ldots, z_N) \right| \leq \frac{3}{N}$$

yielding

$$\mathbb{P}(\sup_{f \in \mathcal{F}} (R(f) - \mathcal{E}_{\mathbb{T}}(f)) \geq C_{\mathbb{T}} + \epsilon) \leq e^{-\frac{2N\epsilon^2}{9}}.$$

### 22.6.3   Rademacher complexity

We now extend the previous definition by computing discrepancies over random two-set partitions of the training set, which have equal size in average. This leads to the empirical Rademacher complexity of the function class. Let $\xi_1, \ldots, \xi_N$ be a sequence of Rademacher random variables (equal to -1 and +1 with equal probability 1/2). Then, the (empirical) Rademacher complexity of the training set $T$ for the model class $\mathcal{F}$ is

$$\mathrm{rad}(T) = \mathbb{E}\left( \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{k=1}^{N} \xi_k r(Y_k, f(X_k)) \,\Big|\, \mathbb{T} = T \right).$$

The mean Rademacher complexity is then the expectation of this quantity over the training set distribution. The Rademacher complexity can be computed with a—costly—Monte-Carlo simulation, in which the best estimator is computed with randomly flipped labels corresponding to the values of $k$ such that $\xi_k = -1$.

This measure of complexity was introduced to the machine learning framework in Koltchinskii and Panchenko [109], Bartlett and Mendelson [19], and Rademacher sums have been extensively studied in relation to empirical processes (cf. Ledoux and Talagrand [117], chapter 4).

One can bound the Rademacher complexity in terms of VC dimension.

**Proposition 22.40** *Let $\mathcal{F}$ be a function class such that $D = VC\text{-}dim(\mathcal{F}) < \infty$. Then*

$$\mathrm{rad}(T) \leq \frac{3}{\sqrt{N}} \sqrt{2D \log(eN/D)}.$$

PROOF  One has, using Hoeffding's inequality

$$\mathbb{P}\left( \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{k=1}^{N} \xi_k r(y_k, f_k) > t \right) \leq |\mathcal{F}(T)| \sup_{f \in \mathcal{F}} \mathbb{P}\left( \frac{1}{N} \sum_{k=1}^{N} \xi_k r(y_k, f_k) > t \right) \leq |\mathcal{F}(T)| e^{-Nt^2/2}.$$

This implies that

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{N} \sum_{k=1}^{N} \xi_k r(y_k, f_k) \right| > t \right) \leq 2|\mathcal{F}(T)| e^{-Nt^2/2}$$

and proposition 22.4 implies

$$\mathrm{rad}(T) \leq \frac{3\sqrt{2|\mathcal{F}(T)|}}{\sqrt{N}}.$$

Therefore if $D = VC\text{-}dim(\mathcal{F}) < \infty$, proposition 22.20 implies

$$\mathrm{rad}(T) \leq \frac{3}{\sqrt{N}} \sqrt{2D \log(eN/D)}.$$

■

We now discuss generalization bounds using Rademacher's complexity. While we still consider binary classification problems (with $\mathcal{R}_Y = \{-1, 1\}$), we will assume that $\mathcal{F}$ contains functions that can take arbitrary scalar values, and the 0–1 loss function becomes $r(y, y') = \mathbf{1}_{yy' \leq 0}$ with $y \in \{-1, 1\}$ and $y' \in \mathbb{R}$. We will also consider functions that dominate this loss, i.e., functions $\rho : \mathcal{R}_Y \times \mathcal{R} \to [0, 1]$ such that

$$r(y, y') \leq \rho(y, y')$$

for all $y \in \mathcal{R}_Y$, $y' \in \mathbb{R}$. Some examples are the margin loss $\rho_h^*(y, y') = \mathbf{1}_{yy' \leq h}$ for $h \geq 0$, or the piecewise linear function

$$\rho_h(y, y') = \begin{cases} 1 & \text{if } yy' \leq 0 \\ 1 - yy'/h & \text{if } 0 \leq yy' \leq h \\ 0 & \text{if } yy' \geq h \end{cases}$$

If $\mathcal{G}$ is a class of functions $g : \mathcal{Z} \to \mathbb{R}$, we will denote

$$Rad_{\mathcal{G}}(z_1, \ldots, z_N) = \frac{1}{N} E\left( \sup_{g \in \mathcal{G}} \sum_{k=1}^{N} \xi_k g(z_k) \right)$$

and

$$\overline{Rad}_{\mathcal{G}}(N) = \mathbb{E}\left( Rad_{\mathcal{G}}(Z_1, \ldots, Z_N) \right).$$

Our previous notation can then be rewritten as $\mathrm{rad}(T) = Rad_{\mathcal{G}}(z_1, \ldots, z_n)$ where $z_i = (x_i, y_i)$ and $\mathcal{G}$ is the space of functions: $g : (x, y) \mapsto r(y, f(x))$ for $f \in \mathcal{F}$. The following theorem is proved in Koltchinskii and Panchenko [109], Bartlett and Mendelson [19].

**Theorem 22.41** *Let $\rho$ be a function dominating the risk function $r(y, y') = \mathbf{1}_{yy' \leq 0}$. Let*

$$\mathcal{G}^\rho = \{(x, y) \mapsto \rho(y, f(x)) - 1 : f \in \mathcal{F}\}$$

*and*

$$\mathcal{E}_T^\rho(f) = \frac{1}{N} \sum_{k=1}^N \rho(y_k, f(x_k)).$$

*Then*

$$\mathbb{P}(\sup_{f \in \mathcal{F}} R(f) \geq \mathcal{E}_T^\rho(f) + 2\overline{Rad}_{\mathcal{G}^\rho}(N) + t) \leq e^{-Nt^2/2}$$

PROOF  For $f \in \mathcal{F}$, we have

$$R(f) - \mathcal{E}^\rho(f) \leq \mathbb{E}(\rho(Y, f(X))) - \mathcal{E}^\rho(f) \leq \Phi(Z_1, \dots, Z_N)$$

where

$$\Phi(Z_1, \dots, Z_N) = \sup_{g \in \mathcal{G}^\rho}\left(\mathbb{E}(g(Z)) - \frac{1}{N}\sum_{k=1}^N g(Z_k)\right).$$

Since changing one variable among $Z_1, \dots, Z_N$ changes $\Phi$ by at most $2/N$, McDiarmid's inequality implies that

$$\mathbb{P}(\Phi(Z_1, \dots, Z_N) - \mathbb{E}(\Phi(Z_1, \dots, Z_N)) \geq t) \leq e^{-Nt^2/2}.$$

Now we have

$$\mathbb{E}(\Phi(Z_1, \dots, Z_N)) = \mathbb{E}\left(\sup_{g \in \mathcal{G}^\rho} \mathbb{E}\left(\frac{1}{N}\sum_{k=1}^N g(Z_k') - \frac{1}{N}\sum_{k=1}^N g(Z_k)\Big| Z_1, \dots, Z_N\right)\right)$$

$$\leq \mathbb{E}\left(\sup_{g \in \mathcal{G}^\rho}\left(\frac{1}{N}\sum_{k=1}^N g(Z_k') - \frac{1}{N}\sum_{k=1}^N g(Z_k)\right)\right)$$

$$\leq \mathbb{E}\left(\mathbb{E}\left(\sup_{g \in \mathcal{G}^\rho}\left(\frac{1}{N}\sum_{k=1}^N \xi_k(g(Z_k') - g(Z_k))\right)\Big| Z, Z'\right)\right)$$

$$\leq 2\mathbb{E}\left(\mathbb{E}\left(\sup_{g \in \mathcal{G}^\rho}\left(\frac{1}{N}\sum_{k=1}^N \xi_k g(Z_k)\right)\Big| Z\right)\right)$$

$$\leq 2\overline{Rad}_{\mathcal{G}^\rho}(N),$$

of which the statement of the theorem is a direct consequence.                                    ∎

### 22.6.4 Algorithmic Stability

Another result using McDiarmid's inequality is proved in Bousquet and Elisseeff [39], and is based on the stability of a classifier when one removes a single example from the training set. As before, we consider training sets $\mathbb{T}$ of size $N$, where $\mathbb{T}$ is a random variable.

For $k \in \{1, \ldots, N\}$, and a training set $T = (x_1, y_1, \ldots, x_N, y_N)$, we let $T^{(k)}$ be the training set with sample $(x_k, y_k)$ removed. One says that the predictor $(T \mapsto \hat{f}_T)$ has uniform stability $\beta_N$ for the loss function $r$ if, for all $T$ of size $N$, all $k \in \{1, \ldots, N\}$, and all $x, y$:

$$|r(\hat{f}_T(x), y) - r(\hat{f}_{T^{(k)}}(x), y)| \le \beta_N. \tag{22.45}$$

With this definition, the following theorem holds.

**Theorem 22.42 (Bousquet and Elisseeff [39])** *Assume that $\hat{f}_T$ has uniform stability $\beta_N$ for training sets of size $N$ and that the loss function $r(Y, f(X))$ is almost surely bounded by $M > 0$. Then, for all $\epsilon > 2\beta_N$, one has*

$$\mathbb{P}(R(\hat{f}_{\mathbb{T}}) \ge \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}) + \epsilon) \le e^{-2N\left(\frac{\epsilon - 2\beta_N}{4N\beta_N + M}\right)^2}.$$

Of course, this theorem is interesting only when $\beta_N$ is small as a function of $N$, i.e., when $N\beta_N$ is bounded.

PROOF Let $Z_i = (X_i, Y_i)$ and $F(Z_1, \ldots, Z_N) = R(\hat{f}_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}})$. We want to apply McDiarmid inequality (theorem 22.13) to $F$, and therefore estimate

$$\delta_k(F) \overset{\Delta}{=} \max_{z_1, \ldots, z_N, z_k'} \left| F(z_1, \ldots, z_N) - F(z_1, \ldots, z_{k-1}, z_k', z_{k+1}, \ldots, z_N) \right|.$$

Introduce a training set $\tilde{T}_k$ in which the variable $z_k$ is replaced by $z_k' = (x_k', y_k')$. Because $\tilde{T}_k^{(k)} = T^{(k)}$, we have

$$
\begin{aligned}
|R(\hat{f}_T) - R(\hat{f}_{\tilde{T}_k})| &\le E(|r(Y, \hat{f}_T(X)) - r(Y, \hat{f}_{\tilde{T}_k}(X)))| \\
&\le E(|r(Y, \hat{f}_T(X)) - r(Y, \hat{f}_{T^{(k)}}(X))|) \\
&\quad + E(|r(Y, \hat{f}_{\tilde{T}_k}(X)) - E(r(Y, \hat{f}_{\tilde{T}_k^{(k)}}(X)))|) \\
&\le 2\beta_N
\end{aligned}
$$

Similarly, we have

$$
\begin{aligned}
|\mathcal{E}_T(\hat{f}_T) - \mathcal{E}_{\tilde{T}_k})(\hat{f}_{\tilde{T}_k})| &\leq \frac{1}{N} \sum_{l \neq k} |r(y_l, \hat{f}_T(x_l),) - r(y_l, \hat{f}_{\tilde{T}_k}(x_l))| \\
&\quad + \frac{1}{N} |r(y_k, \hat{f}_T(x_k)) - r(y_k', \hat{f}_{\tilde{T}_k}(x_k'))| \\
&\leq \frac{1}{N} \sum_{l \neq k} |r(y_l, \hat{f}_T(x_l)) - r(y_l, \hat{f}_{T^{(k)}}(x_l))| \\
&\quad + \frac{1}{N} \sum_{l \neq k} |r(y_l, \hat{f}_{\tilde{T}_k}(x_l)) - r(y_l, \hat{f}_{\tilde{T}_k^{(k)}}(x_l))| + \frac{M}{N} \\
&\leq 2\beta_N + \frac{M}{N}
\end{aligned}
$$

Collecting these results, we find that $\delta_k(F) \leq 4\beta_N + \frac{M}{N}$, so that, by theorem 22.13,

$$
\mathbb{P}\left( R(\hat{f}_{\mathbb{T}}) \geq \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}) + \mathbb{E}(R(\hat{f}_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}))) + \epsilon \right) \leq \exp\left( -\frac{2N\epsilon^2}{(4N\beta_N + M)^2} \right).
$$

It remains to evaluate the expectation in this formula. Introducing as above variables $Z_1', \ldots, Z_N'$ and using the same notation for $\tilde{T}_k$, we have

$$
\mathbb{E}(R(\hat{f}_{\mathbb{T}})) = \mathbb{E}(r(Y_k', f_{\mathbb{T}}(X_k'))) = \mathbb{E}(r(Y_k, f_{\tilde{\mathbb{T}}_k}(X_k))).
$$

Using this, we have

$$
\begin{aligned}
\mathbb{E}(R(\hat{f}_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}})) &= \frac{1}{N} \sum_{k=1}^{N} \mathbb{E}(r(Y_k, f_{\tilde{\mathbb{T}}_k}(X_k)) - r(Y_k, f_{\mathbb{T}}(X_k))) \\
&= \frac{1}{N} \sum_{k=1}^{N} \mathbb{E}(r(Y_k, f_{\tilde{\mathbb{T}}_k}(X_k)) - r(Y_k, f_{\tilde{\mathbb{T}}_k^{(k)}}(X_k))) \\
&\quad + \frac{1}{N} \sum_{k=1}^{N} \mathbb{E}(r(Y_k, f_{\mathbb{T}_k^{(k)}}(X_k)) - r(Y_k, f_{\mathbb{T}}(X_k)))
\end{aligned}
$$

from which one deduces that

$$
|\mathbb{E}(R(\hat{f}_{\mathbb{T}}) - \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}))| \leq 2\beta_N.
$$

We therefore obtain

$$
\mathbb{P}\left( R(\hat{f}_{\mathbb{T}}) \geq \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}) + \epsilon + 2\beta_N \right) \leq \exp\left( -\frac{2N\epsilon^2}{(4N\beta_N + M)^2} \right).
$$

as required.                                                                      ∎

### 22.6.5 PAC-Bayesian bounds

Our final discussion of concentration bounds for the empirical error uses a slightly different paradigm from that discussed so far. The main difference is that, instead of computing one predictor $\hat{f}_T$ from a training set $T$, it would return a random variable with values in $\mathcal{F}$, or, equivalently, a probability distribution on $\mathcal{F}$ (therefore assuming that this space is measurable) that we will denote $\hat{\mu}_T$. The training set error is now defined by:

$$\bar{\mathcal{E}}_T(\mu) = \int \mathcal{E}_T(f) d\mu(f),$$

for any probability distribution $\mu$ on $\mathcal{F}$, while the generalization error is:

$$\bar{R}(\mu) = \int_{\mathcal{F}} R(f) d\mu(f).$$

Our goal is to obtain upper bounds on $\bar{R}(\mu_{\mathbb{T}}) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu_{\mathbb{T}})$ that hold with high probability. In this framework, we have the following result, in which $\mathcal{Q}$ denotes the space of probability distributions on $\mathcal{F}$.

Assume that the loss function $r$ takes its values in $[0,1]$. Recall that $KL(\mu\|\pi)$ is the Kullback-Leibler divergence from $\mu$ to $\pi$, defined by

$$KL(\mu\|\pi) = \int_{\mathcal{F}} \log(\varphi(f))\varphi(f) d\pi(f)$$

if $\mu$ has a density $\varphi$ with respect to $\pi$ and $+\infty$ otherwise. Then, the following theorem holds.

**Theorem 22.43 (McAllester [128])** *With the notation above, for any fixed probability distribution $\pi \in \mathcal{Q}$,*

$$\mathbb{P}\left(\sup_{\mu \in \mathcal{Q}}(\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu)) > \sqrt{t + \frac{KL(\mu\|\pi)}{2N}}\right) \leq 2Ne^{-Nt}. \tag{22.46}$$

Taking $t = \log(2N/\delta)/2N$, the theorem is equivalent to the statement that, with probability $1 - \delta$, one has

$$\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu) \leq \sqrt{\frac{\log 2N/\delta + KL(\mu\|\pi)}{2N}}. \tag{22.47}$$

PROOF We first show that, for any probability distributions $\pi, \mu$ on $\mathcal{F}$, and any function $H$ on $\mathcal{F}$,

$$\int_{\mathcal{F}} H(f) d\mu - \log \int_{\mathcal{F}} e^{H(f)} d\pi \leq KL(\mu\|\pi).$$

Indeed, assume that $\mu$ has a density $\varphi$ with respect to $\pi$ (otherwise the upper bound is infinite) and let

$$\varphi_H = \frac{e^H}{\int_{\mathcal{F}} e^{H(f)} d\pi}.$$

Then

$$KL(\mu\|\pi) - \int_{\mathcal{F}} H(f)d\mu + \log \int_{\mathcal{F}} e^{H(f)}d\pi = \int_{\mathcal{F}} \varphi \log \varphi \, d\pi - \int_{\mathcal{F}} \varphi \log \varphi_H d\pi$$

$$= \int_{\mathcal{F}} \frac{\varphi}{\varphi_H} \left( \log \frac{\varphi}{\varphi_H} \right) \varphi_H d\pi$$

$$= KL(\mu\|\varphi_H \pi) \geq 0,$$

which proves the result (and also shows that one can only have equality when $\varphi = \varphi_H$ $\pi$-almost surely.)

Let $\chi(u) = \max(u, 0)^2$. We can use this inequality to show that, for any probability $Q \in \mathcal{Q}$ and $\lambda > 0$,

$$\lambda \chi(\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu)) \leq \lambda \int_{\mathcal{F}} \chi(R(f) - \mathcal{E}_{\mathbb{T}}(f))d\mu(f) \leq KL(\mu\|\pi) + \log \int_{\mathcal{F}} e^{\lambda \chi(R(f) - \mathcal{E}_{\mathbb{T}}(f))}d\pi$$

where we have applied Jensen's inequality to the convex function $\chi$. This yields

$$e^{\lambda \chi(\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(Q))} \leq e^{KL(\mu\|\pi)} \int_{\mathcal{F}} e^{\lambda \chi(R(f) - \mathcal{E}_{\mathbb{T}}(f))}d\pi.$$

Hoeffding's inequality implies that, for all $f \in \mathcal{F}$ and $t \geq 0$

$$\mathbb{P}(\chi(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > t) = \mathbb{P}(R(f) - \mathcal{E}_{\mathbb{T}}(f) > \sqrt{t}) \leq e^{-2Nt}$$

so that

$$\mathbb{E}\left( e^{\lambda \chi(R(f) - \mathcal{E}_{\mathbb{T}}(f))} \right) = \int_0^\infty \mathbb{P}(\lambda \chi(R(f) - \mathcal{E}_{\mathbb{T}}(f)) > \log t)dt$$

$$\leq 1 + \int_1^{e^\lambda} e^{-\frac{2N \log t}{\lambda}} dt$$

$$= 1 + \int_0^\lambda e^{-\frac{2Nu}{\lambda} + u} du$$

$$= 1 + \lambda \frac{e^{\lambda - 2N} - 1}{\lambda - 2N}.$$

From this and Markov's inequality, we get, for any $\lambda > 0$:

$$\mathbb{P}(\sup_{\mu \in \mathcal{Q}} \chi(\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu)) > t + KL(\mu\|\pi)/\lambda) \leq e^{-\lambda t}\left( 1 + \lambda \frac{e^{\lambda - 2N} - 1}{\lambda - 2N} \right).$$

Taking $\lambda = 2N$ yields

$$\mathbb{P}(\sup_{\mu \in \mathcal{Q}} \chi(\bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu)) > t + KL(\mu\|\pi)/2N) \leq 2Ne^{-2Nt},$$

which implies

$$\mathbb{P}\left(\sup_{\mu \in \mathcal{Q}} \bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu) > \sqrt{t + KL(\mu\|\pi)/2N}\right) \leq 2Ne^{-2Nt},$$

concluding the proof. ∎

**Remark 22.44** Note that the proof, which follows that given in Audibert and Bousquet [15], provides a family of inequalities obtained by taking $\lambda = 2N/c$ in the final step, with $c > 1$. In this case

$$1 + \lambda \frac{e^{\lambda - 2N} - 1}{\lambda - 2N} \leq 1 + \frac{\lambda}{2N - \lambda} = \frac{c}{c - 1}$$

and one gets

$$\mathbb{P}\left(\sup_{\mu \in \mathcal{Q}} \bar{R}(\mu) - \bar{\mathcal{E}}_{\mathbb{T}}(\mu) > \sqrt{t + cKL(\mu\|\pi)/2N}\right) \leq \frac{c}{c - 1} 2Ne^{-2Nt}.$$

◆

**Remark 22.45** One special case of theorem 22.43 is when $\pi$ is a discrete probability measure supported by a subset $\mathcal{F}_0$ of $\mathcal{F}$ and $\mu$ corresponds to a deterministic predictor optimized over $\mathcal{F}_0$, and is therefore a Dirac measure on some element $f \in \mathcal{F}_0$. Because $\delta_f$ has density $\varphi(g) = 1/\pi(g)$ if $g = f$ and 0 otherwise with respect to $\pi$, we have $KL(\delta_f\|\pi) = -\log \pi(f)$ and theorem 22.43 implies that, with probability larger than $1 - \delta$,

$$R(f) - \mathcal{E}_{\mathbb{T}}(f) \leq \sqrt{\frac{\log 2N/\delta - \log \pi(f)}{2N}}.$$

The term $\log 2N$ is however superfluous in this simple context, because one can write, for any $t > 0$

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}_0} R(f) - \mathcal{E}_{\mathbb{T}}(f) \geq \sqrt{t - \frac{\log(\pi(f))}{2N}}\right) \leq \sum_{f \in \mathcal{F}_0} e^{-2N(t\frac{\log(\pi(f))}{2N})} = e^{-2Nt}$$

so that, with probability $1 - \delta$ (letting $t = \log(1/\delta)/2N$), for all $f \in \mathcal{F}_0$:

$$R(f) - \mathcal{E}_{\mathbb{T}}(f) \leq \sqrt{\frac{-\log \delta - \log \pi(f)}{2N}}.$$

◆

## 22.7    Application to model selection

We now describe how the previous results can, in principle, be applied to model selection [20]. We assume that we have a countable family of nested models classes $(\mathcal{F}^{(j)}, j \in \mathcal{J})$. Denote, as usual, by $\mathcal{E}_T(f)$ the empirical prediction error in the training set for a given function $f$. We will denote by $\hat{f}_T^{(j)}$ a minimizer of the in-sample error for $\mathcal{F}^{(j)}$, such that

$$\mathcal{E}_T(\hat{f}_T^{(j)}) = \min_{f \in \mathcal{F}^{(j)}} \mathcal{E}_T(f).$$

In the model selection problem, one would like to determine the best model class, $j = j(T)$, such that the prediction error $R(\hat{f}_T^{(j)})$ is minimal, or, more realistically, determine $j^*$ such that $R(\hat{f}_T^{(j_*)})$ is not too far from the optimal one.

We will consider penalty-based methods in which one minimizes $\tilde{\mathcal{E}}_T(f) = \mathcal{E}_T(f) + C_T(j)$ to determine $j(T)$. The penalty, $C_T$, may also be data-dependent, and will therefore be a random variable. The previous concentration inequalities provided highly probable upper-bounds for $R(\hat{f}_{\mathbb{T}}^{(j)})$, each exhibiting a random variable $\Gamma_{\mathbb{T}}^{(j)}$ that is larger than $R(\hat{f}_{\mathbb{T}}^{(j)})$ with probability close to one. More precisely, we obtained inequalities taking the form (when applied to $\mathcal{F}^{(j)}$)

$$\mathbb{P}(R_{\mathbb{T}}(\hat{f}^{(j)}) \geq \Gamma_{\mathbb{T}}^{(j)} + t) \leq c_j e^{-mt^2} \tag{22.48}$$

for some known constants $c_j$ and $m$. For example, the VC-dimension bounds have $\Gamma_{\mathbb{T}}^{(j)} = \mathcal{E}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}^{(j)})$, $c_j = 2\mathcal{S}_{\mathcal{F}^{(j)}}(2N)$ and $m = N/8$.

Given such inequalities, one can develop a model selection strategy that relies on *a priori* weights, provided by a sequence $\pi_j$ of positive numbers such that $\sum_{j \in \mathcal{J}} \pi_j = 1$. Define

$$\tilde{\pi}_j = \frac{\pi_j/c_j}{\sum_{j'=1}^{\infty} \pi_{j'}/c_{j'}},$$

and let

$$C_T^{(j)} = \Gamma_T^{(j)} - \mathcal{E}_T(\hat{f}_T^{(j)}) + \sqrt{-\frac{\log \tilde{\pi}_j}{m}}$$

yielding a penalty-based method that requires the minimization of

$$\tilde{\mathcal{E}}_T(f) = (\mathcal{E}_T(f) - \mathcal{E}_T(\hat{f}_T^{(j)})) + \Gamma_T^{(j)} + \sqrt{-\frac{\log \tilde{\pi}_j}{m}}.$$

The selected model class is then $\mathcal{F}^{(j^*)}$ where $j^*$ minimizes $\Gamma_T^{(j)} + \sqrt{-\frac{\log \tilde{\pi}_j}{2m}}$.

The same proof as that provided at the end of section 22.6.5 justifies this procedure. Indeed, for $t > 0$,

$$
\mathbb{P}\left(R(\hat{f}_{\mathbb{T}}) - \tilde{\mathcal{E}}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}) \geq t\right) \leq \mathbb{P}\left(\max_j (R(\hat{f}_{\mathbb{T}}^{(j)}) - \tilde{\mathcal{E}}_{\mathbb{T}}(\hat{f}_{\mathbb{T}}^{(j)})) \geq t\right)
$$

$$
\leq \mathbb{P}\left(\max_j (R(\hat{f}_{\mathbb{T}}^{(j)}) \geq R_j^* + t + \sqrt{-\frac{\log \tilde{\pi}_j}{m}}\right)
$$

$$
\leq \tilde{c} \sum_j \pi_j e^{-mt^2}
$$

$$
\leq \tilde{c} e^{-mt^2}
$$

with $\tilde{c} = \sum_{j=1}^{\infty} \pi_j / c_j$.

# Bibliography

[1] Pierre-Antoine Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.

[2] Hirotugu Akaike. Information theory and an extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory, 1973*. Akademiai Kaido, 1973.

[3] Stéphanie Allassonniere and Laurent Younes. A stochastic algorithm for probabilistic independent component analysis. *The Annals of Applied Statistics*, 6 (1):125–160, 2012.

[4] Noga Alon, Shai Ben-David, Nicolo Cesa-Bianchi, and David Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *Journal of the ACM (JACM)*, 44(4):615–631, 1997.

[5] Mauricio A. Álvarez, Lorenzo Rosasco, and Neil D. Lawrence. Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3):195–266, 2012. ISSN 1935-8237. doi: 10.1561/2200000036.

[6] Yali Amit. Convergence properties of the gibbs sampler for perturbations of gaussians. *The Annals of Statistics*, 24(1):122–140, 1996.

[7] Yali Amit and Donald Geman. Shape quantization and recognition with randomized trees. *Neural computation*, 9(7):1545–1588, 1997.

[8] Alano Ancona, Donald Geman, Nobuyuki Ikeda, and D Geman. Random fields and inverse problems in imaging. In *Ecole d'ete de Probabilites de Saint-Flour XVIII-1988*, pages 115–193. Springer, 1990.

[9] Brian D.O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.

[10] Martin Anthony and Peter L. Bartlett. *Neural network learning: Theoretical foundations*. cambridge university press, 2009.

[11] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pages 214–223. JMLR.org, 2017. event-place: Sydney, NSW, Australia.

[12] Nachman Aronszajn. Theory of Reproducing Kernels. *Trans. Am. Math. Soc.*, 68:337–404, 1950.

[13] Krishna B. Athreya, Hani Doss, and Jayaram Sethuraman. On the convergence of the markov chain simulation method. *The Annals of Statistics*, 24(1):69–100, 1996.

[14] Hagai Attias. A Variational Baysian Framework for Graphical Models. In *NIPS*, volume 12. Citeseer, 1999.

[15] Jean-Yves Audibert and Olivier Bousquet. Combining pac-bayesian and generic chaining bounds. *Journal of Machine Learning Research*, 8(Apr):863–889, 2007.

[16] Adrian Barbu and Song-Chun Zhu. Generalizing swendsen-wang to sampling arbitrary posterior probabilities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(8):1239–1253, 2005.

[17] Viorel Barbu. *Differential equations*. Springer, 2016.

[18] Peter Bartlett and John Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. *Advances in Kernel methods—support vector learning*, pages 43–54, 1999.

[19] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

[20] Peter L. Bartlett, Stéphane Boucheron, and Gábor Lugosi. Model selection and error estimation. *Machine Learning*, 48:85–113, 2002.

[21] Peter L. Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.

[22] Amir Beck. *Introduction to nonlinear optimization: Theory, algorithms, and applications with MATLAB*. SIAM, 2014.

[23] Michel Benaïm. Dynamics of stochastic approximation algorithms. In *Seminaire de probabilites XXXIII*, pages 1–68. Springer, 1999.

[24] George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.

[25] Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic approximations*, volume 22. Springer Science & Business Media, 2012.

[26] Nils Berglund. Long-time dynamics of stochastic differential equations. *arXiv preprint arXiv:2106.12998*, 2021.

[27] Dimitri Bertsekas. *Convex optimization theory*, volume 1. Athena Scientific, 2009.

[28] Rajendra Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.

[29] Peter J. Bickel and Kjell A. Doksum. *Mathematical statistics: basic ideas and selected topics, volume I*, volume 117. CRC Press, 2015.

[30] Peter J. Bickel, Ya'acov Ritov, and Alexandre B. Tsybakov. Simultaneous analysis of lasso and dantzig selector. 2009.

[31] Patrick Billingsley. *Probability and measure*. John Wiley & Sons, 2008.

[32] Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.

[33] Salomon Bochner. Vorlesungen über fouriersche integrale. *Bull Amer Math Soc*, 39:184, 1933.

[34] Vladimir I. Bogachev. *Measure Theory*. Springer, 2007.

[35] Joseph-Frédéric Bonnans, Jean Charles Gilbert, Claude Lemaréchal, and Claudia A. Sagastizábal. *Numerical optimization: theoretical and practical aspects*. Springer Science & Business Media, 2006.

[36] Ingwer Borg and Patrick J.F. Groenen. *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.

[37] Jonathan Borwein and Adrian S. Lewis. *Convex analysis and nonlinear optimization: theory and examples*. Springer Science & Business Media, 2010.

[38] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. A sharp concentration inequality with applications. *Random Structures & Algorithms*, 16(3):277–292, 2000.

[39] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of machine learning research*, 2(Mar):499–526, 2002.

[40] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.

[41] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.

[42] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.

[43] Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A. Olshen. *Classification and regression trees*. CRC press, 1984.

[44] Dmitri Burago, Iu D. Burago, Yuri Burago, Sergei A. Ivanov, and Sergei Ivanov. *A course in metric geometry*, volume 33. American Mathematical Soc., 2001.

[45] Jian-Feng Cai, Emmanuel J. Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on optimization*, 20(4):1956–1982, 2010.

[46] Tadeusz Caliński and Jerzy Harabasz. A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods*, 3(1):1–27, 1974.

[47] Emmanuel J. Candes and Terence Tao. Decoding by linear programming. *IEEE Trans. information theory*, 51(12):4203–4215, 2005.

[48] Emmanuel J. Candes and Terence Tao. The dantzig selector: statistical estimation when $p$ is much larget. *Annals of statistics*, 35, 2007.

[49] Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):11, 2011.

[50] John Canny. Gap: a factor model for discrete data. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 122–129, 2004.

[51] B. Chalmond. An iterative Gibbsian technique for reconstruction of m-ary images. *Pattern recognition*, 22(6):747–761, 1989. ISSN 0031-3203.

[52] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 6571–6583. Curran Associates, Inc., 2018.

[53] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.

[54] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.

[55] Thomas M. Cover and Joy A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

[56] Robert G. Cowell, A. Philip Dawid, Steffen L. Lauritzen, and David J. Spiegelhalter. *Probabilistic networks and expert systems*. Springer, 2007.

[57] George Darmois. Analyse générale des liaisons stochastiques: etude particulière de l'analyse factorielle linéaire. *Revue de l'Institut international de statistique*, pages 2–8, 1953.

[58] Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the em algorithm. *Annals of statistics*, pages 94–128, 1999.

[59] Amir Dembo and Ofer Zeitouni. Large deviations techniques and applications. 1998. *Applications of Mathematics*, 38, 2011.

[60] Luc Devroye, Lázló Györfi, and Gábor Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.

[61] Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional gaussians with the same mean. *arXiv preprint arXiv:1810.08693*, 2018.

[62] Jean Dieudonné. *Infinitesimal Calculus*. Houghton Mifflin, 1971.

[63] Edsger W. Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1):269–271, 1959. ISSN 0029-599X.

[64] Petros Drineas, Alan Frieze, Ravi Kannan, Santosh Vempala, and Vishwanathan Vinay. Clustering large graphs via the singular value decomposition. *Machine learning*, 56:9–33, 2004.

[65] Simon Duane, Anthony D. Kennedy, Brian J. Pendleton, and Duncan Roweth. Hybrid monte carlo. *Physics letters B*, 195(2):216–222, 1987.

[66] Richard M. Dudley. *Real analysis and probability*. Chapman and Hall/CRC, 2018.

[67] Marie Duflo. *Random iterative models*, volume 34. Springer Science & Business Media, 2013.

[68] HA Eiselt, Carl-Louis Sandblom, et al. *Nonlinear optimization: Methods and applications*. Springer, 2019.

[69] Stewart N. Ethier and Thomas G. Kurtz. Markov processes: Characterization and convergence. 1986.

[70] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

[71] James A. Fill. An interruptible algorithm for perfect sampling via Markov chains. *The Annals of Applied Probability*, 8(1):131–162, 1998.

[72] P. Thomas Fletcher and Sarang Joshi. Principal geodesic analysis on symmetric spaces: Statistics of diffusion tensors. In *Computer vision and mathematical methods in medical and biomedical image analysis*, pages 87–98. Springer, 2004.

[73] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997. Publisher: Elsevier.

[74] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.

[75] Jerome H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001. Publisher: JSTOR.

[76] Dan Geiger and Judea Pearl. On the logic of causal models. In *Machine intelligence and pattern recognition*, volume 9, pages 3–14. Elsevier, 1990.

[77] Dan Geiger, Thomas Verma, and Judea Pearl. Identifying independence in bayesian networks. *Networks*, 20(5):507–534, August 1990. ISSN 00283045. doi: 10.1002/net.3230200504.

[78] Donald Geman, Christian d'Avignon, Daniel Q. Naiman, and Raimond L Winslow. Classifying gene expression profiles from pairwise mrna comparisons. *Statistical applications in genetics and molecular biology*, 3(1):1–19, 2004.

[79] Stuart Geman and Donald Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence*, (6):721–741, 1984.

[80] Stuart Geman and Chii-Ruey Hwang. Nonparametric maximum likelihood estimation by the method of sieves. *The Annals of Statistics*, pages 401–414, 1982.

[81] M. Gondran and M. Minoux. *Graphs and algorithms*. John Wiley & Sons, 1983.

[82] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[83] Ulf Grenander. *Abstract Inference*. Wiley, 1981.

[84] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of Wasserstein GANs. In *Advances in neural information processing systems*, pages 5767–5777, 2017.

[85] Madan M. Gupta and J. Qi. Theory of T-norms and fuzzy inference methods. *Fuzzy Sets and Systems*, 40(3):431–450, April 1991. ISSN 0165-0114.

[86] Te Sun Han. Nonnegative entropy measures of multivariate symmetric correlations. *Information and Control*, 36:133–156, 1978.

[87] Trevor Hastie, Robert Tibshirani, and Jerome H. Friedman. *The elements of statistical learning*. Springer, 2003.

[88] W. Keith Hastings. Monte carlo sampling methods using markov chains and their applications. 1970.

[89] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778. IEEE, 2016.

[90] Geoffrey E. Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15:857–864, 2002.

[91] Leslie M. Hocking. *Optimal Control: An Introduction to the Theory with Applications*. Oxford University Press, 1991.

[92] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. In *The Collected Works of Wassily Hoeffding*, pages 409–426. Springer, 1994.

[93] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge university press, 2012.

[94] Aapo Hyvärinen. New approximations of differential entropy for independent component analysis and projection pursuit. In *Advances in neural information processing systems*, pages 273–279, 1998.

[95] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.

[96] Nobuyuki Ikeda and Shinzo Watanabe. *Stochastic differential equations and diffusion processes*. Elsevier, 1981.

[97] Tommi Sakari Jaakkola. *Variational methods for inference and estimation in graphical models*. PhD Thesis, Massachusetts Institute of Technology, 1997.

[98] Vojtech Jarnik. O jistem problemu minimalnim (about a certain minimal problem). *Prace Moravske Prirodovedecke Spolecnosti*, 6:57–63, 1930.

[99] Finn Jensen and Frank Jensen. Optimal junction trees. In *Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence*, pages 360–366, 1994.

[100] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.

[101] Abram M. Kagan, Calyampudi Radhakrishna Rao, and Yurij Vladimirovich Linnik. Characterization problems in mathematical statistics. 1973.

[102] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[103] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. 2014.

[104] Diederik P. Kingma and Max Welling. An Introduction to Variational Autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. Publisher: Now Publishers, Inc.

[105] John Kingman. Completely random measures. *Pacific Journal of Mathematics*, 21(1):59–78, 1967.

[106] Peter E. Kloeden and Eckhard Platen. *Numerical solutions of stochastic differential equations*. Springer, 1992.

[107] Ivan Kobyzev, Simon J.D. Prince, and Marcus A. Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.

[108] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. The MIT Press, 2009.

[109] Vladimir Koltchinskii and Dmitry Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1):1–50, 2002.

[110] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

[111] Wojtek J. Krzanowski and Y.T. Lai. A criterion for determining the number of groups in a data set using sum-of-squares clustering. *Biometrics*, pages 23–34, 1988.

[112] Estelle Kuhn and Marc Lavielle. Coupling a stochastic approximation version of EM with an MCMC procedure. *ESAIM: Probability and Statistics*, 8:115–131, 2004. Publisher: EDP Sciences.

[113] Harold Kushner and G. George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.

[114] Steffen L Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.

[115] Yann LeCun and Yoshua Bengio. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10): 1995, 1995.

[116] Yann LeCun, Bernhard Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne Hubbard, and Lawrence D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.

[117] Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 1991.

[118] Erich L. Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.

[119] Benedict Leimkuhler and Sebastian Reich. *Simulating Hamiltonian Dynamics*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2005.

[120] Lennart Ljung. Analysis of recursive stochastic algorithms. *IEEE Transactions on Automatic Control*, 22(4):551–575, August 1977. ISSN 1558-2523. Conference Name: IEEE Transactions on Automatic Control.

[121] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.

[122] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

[123] Jack Macki and Aaron Strauss. *Introduction to Optimal Control Theory*. Springer Science & Business Media, 2012.

[124] James MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.

[125] Adam a Margolin, Ilya Nemenman, Katia Basso, Chris Wiggins, Gustavo Stolovitzky, Riccardo Dalla Favera, and Andrea Califano. ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC bioinformatics*, 7 Suppl 1:S7, January 2006. ISSN 1471-2105. doi: 10.1186/1471-2105-7-S1-S7.

[126] Enzo Marinari and Giorgio Parisi. Simulated tempering: a new monte carlo scheme. *Europhysics letters*, 19(6):451, 1992.

[127] Pascal Massart. *Concentration inequalities and model selection*, volume 6. Springer, 2007.

[128] David A. McAllester. Pac-bayesian model averaging. In *COLT*, volume 99, pages 164–170. Citeseer, 1999.

[129] James A. McHugh. *Algorithmic graph theory*. New Jersey: Prentice-Hall Inc, 1990.

[130] Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:1802.03426 [cs, stat]*, September 2020. arXiv: 1802.03426.

[131] Henry P. McKean. *Stochastic integrals*, volume 353. American Mathematical Society, 1969.

[132] Kerrie L. Mengersen and Richard L. Tweedie. Rates of convergence of the hastings and metropolis algorithms. *The annals of Statistics*, 24(1):101–121, 1996.

[133] Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, Augusta H. Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092, 1953.

[134] Sean P. Meyn and Richard L. Tweedie. Stability of markovian processes ii: Continuous-time processes and sampled chains. *Advances in Applied Probability*, 25(3):487–517, 1993.

[135] Sean P. Meyn and Richard L. Tweedie. Stability of markovian processes iii: Foster–lyapunov criteria for continuous-time processes. *Advances in Applied Probability*, 25(3):518–548, 1993.

[136] Sean P. Meyn and Richard L. Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.

[137] Leon Mirsky. A trace inequality of john von neumann. *Monatshefte für mathematik*, 79(4):303–306, 1975.

[138] Michel Métivier and Pierre Priouret. Théorèmes de convergence presque sure pour une classe d'algorithmes stochastiques à pas décroissant. *Probability Theory and related fields*, 74(3):403–428, 1987. Publisher: Springer.

[139] Elizbar A. Nadaraya. On estimating regression. *Theory of Probability & Its Applications*, 9(1):141–142, 1964.

[140] Radford M. Neal. Sampling from multimodal distributions using tempered transitions. *Statistics and computing*, 6:353–366, 1996.

[141] Radford M. Neal. Markov chain sampling methods for dirichlet process mixture models. *Journal of computational and graphical statistics*, 9(2):249–265, 2000.

[142] Radford M. Neal. Mcmc using hamiltonian dynamics. *arXiv preprint arXiv:1206.1901*, 2012.

[143] Radford M. Neal and Geoffrey E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer, 1998.

[144] R. E. Neapolitan. *Learning Bayesian networks*. Prentice Hall, 2004.

[145] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust Stochastic Approximation Approach to Stochastic Programming. *SIAM Journal on Optimization*, 19(4):1574–1609, January 2009. ISSN 1052-6234. Publisher: Society for Industrial and Applied Mathematics.

[146] Jorge Nocedal and Stephen J. Wright. *Nonlinear Equations*. Springer, 2006.

[147] Esa Nummelin. *General irreducible Markov chains and non-negative operators*. Number 83. Cambridge University Press, 2004.

[148] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(57):1–64, 2021.

[149] Panos M. Pardalos and Jue Xue. The maximum clique problem. *Journal of Global Optimization*, 4(3):301–328, 1994. ISSN 0925-5001.

[150] Emanuel Parzen. On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33(3):1065–1076, 1962.

[151] Judea Pearl. *Probabilistic reasoning in intelligent systems*. Morgan Kaufmann, 1988, 2012.

[152] Jiming Peng and Yu Wei. Approximating k-means-type clustering via semidefinite programming. *SIAM journal on optimization*, 18(1):186–205, 2007.

[153] Jiming Peng and Yu Xia. A new theoretical framework for k-means-type clustering. *Foundations and advances in data mining*, pages 79–96, 2005. Publisher: Springer.

[154] Odile Pons. *Functional estimation for density, regression models and processes*. World scientific, 2011.

[155] Robert C. Prim. Shortest connection networks and some generalizations. *Bell system technical journal*, 36(6):1389–1401, 1957.

[156] James G. Propp and David B. Wilson. Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures and Algorithms*, 9(1&2):223–252, 1996.

[157] James G. Propp and David B. Wilson. How to get a perfectly random sample from a generic Markov chain and generate a random spanning tree of a directed graph. *Journal of Algorithms*, 27:170–217, 1998.

[158] Jim O. Ramsay and Bernard W. Silverman. *Functional Data Analysis*. Springer-Verlag, 1997.

[159] BLS Prakasa Rao. *Nonparametric functional estimation*. Academic press, 1983.

[160] Daniel Revuz. *Markov chains*. Elsevier, 2008.

[161] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.

[162] Jorma Rissanen. *Stochastic complexity in statistical inquiry*. World Scientific, 1989.

[163] Herbert Robbins and Sutton Monro. A stochastic approximation method. In *Herbert Robbins Selected Papers*, pages 102–109. Springer, 1985.

[164] Gareth O. Roberts and Nicholas G. Polson. On the geometric convergence of the gibbs sampler. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 56(2):377–384, 1994.

[165] Gareth O. Roberts and Jeffrey S. Rosenthal. General state space markov chains and mcmc algorithms. *Probability Surveys*, 1:20–71, 2004.

[166] Gareth O. Roberts and Richard L. Tweedie. Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli*, 2(4):341–363, 1996.

[167] R. Tyrrell Rockafellar. *Convex analysis*, volume 18. Princeton university press, 1970.

[168] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Lecture Notes in Computer Science, pages 234–241, Cham, 2015. Springer International Publishing. ISBN 978-3-319-24574-4.

[169] Kenneth Rose, Eitan Gurewitz, and Geoffrey Fox. A deterministic annealing approach to clustering. *Pattern Recognition Letters*, 11(9):589–594, 1990.

[170] Peter J. Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.

[171] Walter Rudin. *Real and Complex Analysis*. Tata McGraw Hill, 1966.

[172] Robert E. Schapire. The strength of weak learnability. *Machine learning*, 5(2):197–227, 1990.

[173] Isaac J. Schoenberg. Metric spaces and completely monotone functions. *Annals of Mathematics*, pages 811–841, 1938.

[174] Gideon Schwarz. Estimating the dimension of a model. *The annals of statistics*, 6(2):461–464, 1978.

[175] Claude E. Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.

[176] Claude E. Shannon. Communication in the presence of noise. *Proc. Institute of Radio Engineers*, 37(1):10–21, 1949.

[177] Simon J. Sheather and Michael C. Jones. A reliable data-based bandwidth selection method for kernel density estimation. *Journal of the Royal Statistical Society: Series B (Methodological)*, 53(3):683–690, 1991.

[178] Bernard W. Silverman. *Density estimation for statistics and data analysis*. Chapman et Hall, 1998.

[179] Viktor Pavlovich Skitovich. Linear forms of independent random variables and the normal distribution law. *Izvestiya Rossiiskoi Akademii Nauk. Seriya Matematicheskaya*, 18(2):185–200, 1954.

[180] Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR, 2020.

[181] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

[182] Hugo Steinhaus. Sur la division des corp materiels en parties. *Bull. Acad. Polon. Sci*, 1(804):801, 1956.

[183] Charles J. Stone. Consistent nonparametric regression. *The annals of statistics*, pages 595–620, 1977.

[184] Mervyn Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the royal statistical society: Series B (Methodological)*, 36(2): 111–133, 1974.

[185] Catherine A. Sugar and Gareth M. James. Finding the number of clusters in a dataset: An information-theoretic approach. *Journal of the American Statistical Association*, 98(463):750–763, 2003.

[186] Robert H. Swendsen and Jian-Sheng Wang. Nonuniversal critical dynamics in monte carlo simulations. *Physical review letters*, 58(2):86, 1987.

[187] Esteban G. Tabak and Eric Vanden-Eijnden. Density estimation by dual ascent of the log-likelihood. 2010.

[188] Michel Talagrand. *The generic chaining: upper and lower bounds of stochastic processes*. Springer Science & Business Media, 2006.

[189] Michel Talagrand. *Upper and lower bounds for stochastic processes: modern methods and classical problems*, volume 60. Springer Science & Business Media, 2014.

[190] Aik Choon Tan, Daniel Q. Naiman, Lei Xu, Raimond L. Winslow, and Donald Geman. Simple decision rules for classifying human cancers from gene expression profiles. *Bioinformatics*, 21(20):3896–3904, 2005.

[191] Robert Tibshirani, Guenther Walther, and Trevor Hastie. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2):411–423, 2001.

[192] Luke Tierney. Markov Chains for Exploring Posterior Distributions. *Annals of Statistics*, 22(4):1701–1728, December 1994. ISSN 0090-5364, 2168-8966. Publisher: Institute of Mathematical Statistics.

[193] Laurens Van Der Maaten. Accelerating t-sne using tree-based algorithms. *The Journal of Machine Learning Research*, 15(1):3221–3245, 2014.

[194] Aad W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.

[195] Aad W. Van der Vaart and John A. Wellner. *Weak convergence and empirical processes with applications to statistics*. Springer, 1996.

[196] Vladimir Vapnik. *Statistical learning theory. 1998*. Wiley, New York, 1998.

[197] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.

[198] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[199] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.

[200] Rene Vidal, Yi Ma, and Shankar Sastry. Generalized principal component analysis (gpca). *IEEE transactions on pattern analysis and machine intelligence*, 27(12):1945–1959, 2005.

[201] Grace Wahba. *Spline Models for Observational Data*. SIAM, 1990.

[202] Geoffrey S. Watson. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics*, *Series A*, pages 359–372, 1964.

[203] Gerhard Winkler. *Image analysis*, *random fields and Markov chain Monte Carlo methods*. Springer, 1995,2003.

[204] Stephen J. Wright and Benjamin Recht. *Optimization for data analysis*. Cambridge University Press, 2022.

[205] Kôsaku Yosida. *Functional Analysis*. Springer, 1970.

[206] Laurent Younes. Estimation and annealing for gibbsian fields. *Ann. de l'Inst. Henri Poincaré*, 2, 1988.

[207] Laurent Younes. Parametric inference for imperfectly observed gibbsian fields. *Prob. Thry. Rel. Fields*, 82:625–645, 1989.

[208] Laurent Younes. On the convergence of markovian stochastic algorithms with rapidly decreasing ergodicity rates. *Stochastics: An International Journal of Probability and Stochastic Processes*, 65(3-4):177–228, 1999.

[209] Laurent Younes. Diffeomorphic learning. *Journal of Machine Learning Research*, 21:1 – 28, 2020.

[210] Lotfi A. Zadeh. Fuzzy sets. In *Fuzzy sets*, *fuzzy logic, and fuzzy systems: selected papers by Lotfi A Zadeh*, pages 394–432. World Scientific, 1996.