

---

SB2.2/SM4 Statistical Machine Learning  
Lecture notes  
Part I: Unsupervised Learning

---

François Caron

University of Oxford, Hilary Term 2021

Version of January 17, 2021  
Please report typos to [caron@stats.ox.ac.uk](mailto:caron@stats.ox.ac.uk).

# Aims and Structure of this course

The aims of this course are

- To learn a number of different machine learning methods and gain understanding about their statistical foundations;
- To learn to identify and use the appropriate method for a given dataset and a given task;
- To learn how to use the relevant Python libraries/modules to analyse data, interpret results and evaluate the methods.

The course is structured in two main parts. In the first part, we will cover unsupervised learning (dimensionality reduction, feature extraction and clustering). In the second part of the course, we will cover supervised learning (classification and regression).

## References

These lecture notes are based on the slides from previous lecturers who have taught this course, and on the following books:

- L. Wasserman. All of Statistics. Springer, 2010.
- K. Murphy. Machine Learning. A probabilistic perspective. The MIT Press, 2012
- C. M. Bishop. Pattern Recognition and Machine Learning. Springer, 2006.
- T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical Learning. Springer, 2009.
- I. Goodfellow, Y. Bengio and A. Courville. Deep Learning. The MIT Press, 2016.

## Software

This document includes a number of examples written in the programming language Python. Python is certainly the most popular programming language for machine learning; it can be freely downloaded and installed. If you are new to Python, the following website contains useful information for beginners: <https://www.python.org/about/gettingstarted/>.

A number of Python code editors are freely available: PyCharm, Spyder, Atom, Jupyter. You may want to consider installing Anaconda, which is a distribution of Python and R programming languages for data science and machine learning. It comes with code editors such as Spyder and Jupyter for Python.

Python contains a very rich collection of libraries/packages for data science/machine learning. We will use the following libraries in this course:

- Numpy and Scipy: Core toolboxes for scientific computing (mathematical functions, linear algebra routines, random number generators, optimisation, etc.)
- Matplotlib: Visualisation
- Pandas: Data manipulation
- Seaborn: Statistical visualisation
- Scikitlearn: Machine Learning
- Pytorch: Deep Learning

These libraries can be loaded with the following commands.

```
import numpy as np
import scipy as sp
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import sklearn as skl
import torch
```

# Background Material

An excellent resource is the matrix cookbook<sup>1</sup>. Let  $x_1, \dots, x_p$  be scalar values. A  $p$ -dimensional row vector is noted as  $(x_1, \dots, x_p)$ . A  $p$ -dimensional column vector is noted as  $(x_1, \dots, x_p)^\top$ . By default, a vector means a column vector. The L2 norm of a vector  $x = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$  is defined as

$$\|x\| = \sqrt{x^\top x} = \sqrt{\sum_{j=1}^p x_j^2}$$

## 0.1 Derivatives

Let  $J : \mathbb{R}^p \rightarrow \mathbb{R}$  be a scalar function. For a vector  $x = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ , the vector of partial derivatives, or gradient, is the  $p$ -dimensional vector given by

$$\nabla_x J(x) = \frac{\partial J(x)}{\partial x} = \begin{pmatrix} \frac{\partial J(x)}{\partial x_1} \\ \frac{\partial J(x)}{\partial x_2} \\ \vdots \\ \frac{\partial J(x)}{\partial x_p} \end{pmatrix}$$

The only thing you need to remember/know is that  
 $\frac{\partial(x^\top y)}{\partial x} = y$   
and the chain rule, which goes as

$$\underbrace{\frac{d(f(x, y))}{dx}}_{\text{Hence,}} = \frac{\partial(f(x, y))}{\partial x} + \frac{d(x^\top x)}{dx} \frac{\partial(f(x, y))}{\partial y}$$

$$\frac{d(b^\top x)}{dx} = \frac{d(x^\top b)}{dx} = b$$

$$\frac{d(x^\top Ax)}{dx} = \frac{\partial(x^\top y)}{\partial x} + \frac{d(y(x)^\top)}{dx} \frac{\partial(x^\top y)}{\partial y}$$

where  $y = Ax$ . And then, that is,

$$\frac{d(x^\top Ax)}{dx} = \frac{\partial(x^\top y)}{\partial x} + \frac{d(y(x)^\top)}{dx} \frac{\partial(x^\top y)}{\partial y} = y + \frac{d(x^\top A^\top)}{dx} x = y + A^\top x = (A + A^\top)x$$

## 0.2 Properties of the trace

Let  $A = (A_{ij})$  be a real-valued  $p$ -by- $p$  matrix. The trace of  $A$  is defined as

$$\text{trace}(A) = \sum_{i=1}^p A_{ii} = \sum_{i=1}^p \lambda_i \quad (4)$$

where  $\lambda_1, \dots, \lambda_p$  are the eigenvalues of  $A$ . For two matrices  $A$  and  $B$ ,

$$\text{trace}(A) = \text{trace}(A^\top) \quad (5)$$

$$\text{trace}(AB) = \text{trace}(BA) \quad (6)$$

$$\text{trace}(A + B) = \text{trace}(A) + \text{trace}(B) \quad (7)$$

## 0.3 Random vectors

Let  $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$  be a random vector. We denote  $\mathbb{E}[X]$  the mean of  $X$  and

$$\text{cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top]$$

the  $p$ -by- $p$  covariance matrix of  $X$ . Let  $\Sigma = \text{cov}(X)$ . The total variance of the vector  $X$  is defined as

$$\text{TV}(X) = \mathbb{E}(\|X - \mathbb{E}[X]\|^2) = \sum_{j=1}^p \text{var}(X_j) = \text{trace}(\Sigma).$$

<sup>1</sup><https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

# 1 — Introduction

## Week 1.

The objective of unsupervised learning is to find a latent representation of the data that captures relevant information. This representation can be used to find informative ways to represent graphically the data, discover subgroups amongst the data, or reduce the dimensionality.

### 1.1 Notations

The data will consist of  $p$  variables (features/attributes/dimensions) on  $n$  examples (inputs/observations). Let  $\mathbf{X} = (x_{ij})$  be a  $n \times p$  matrix with  $x_{ij}$  := the  $j$ -th variable for the  $i$ -th example

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix}. \quad \mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix}^T$$

Denote the  $i$ -th data item by  $x_i \in \mathbb{R}^p$  (we will treat it as a column vector: it is the transpose of the  $i$ -th row of  $\mathbf{X}$ ). Assume  $x_1, \dots, x_n$  are independently and identically distributed samples of a random vector  $X$  over  $\mathbb{R}^p$ . The  $j$ -th dimension of  $X$  will be denoted  $X_j$ .

**Exploratory data analysis.** Exploratory data analysis is a useful tool to look at data summaries. As an illustration, consider the following dataset. Campbell (1974) studied rock crabs of the genus *leptograpsus*. One species, *L. variegatus*, had been split into two new species according to their colour: orange and blue. Preserved specimens lose their colour, so it was hoped that morphological differences would enable museum material to be classified. Data are available on 50 specimens of each sex of each species. Each specimen has measurements on:

- the width of the frontal lobe FL,
- the rear width RW,
- the length along the carapace midline CL,
- the maximum width CW of the carapace, and
- the body depth BD in mm.

in addition to colour/species and sex (we will later view these as labels, but will ignore for now). Here is a subset of the data:

	FL	RW	CL	CW	BD
1	8.1	6.7	16.1	19.0	7.0
2	8.8	7.7	18.1	20.8	7.4
3	9.2	7.8	19.0	22.4	7.7
4	9.6	7.9	20.1	23.1	8.2
5	9.8	8.0	20.3	23.0	8.2
6	10.8	9.0	23.0	26.5	9.8
7	11.1	9.9	23.8	27.1	9.8
8	11.6	9.1	24.5	28.4	10.4
9	11.8	9.6	24.2	27.8	9.7
10	11.8	10.5	25.2	29.3	10.3
...					

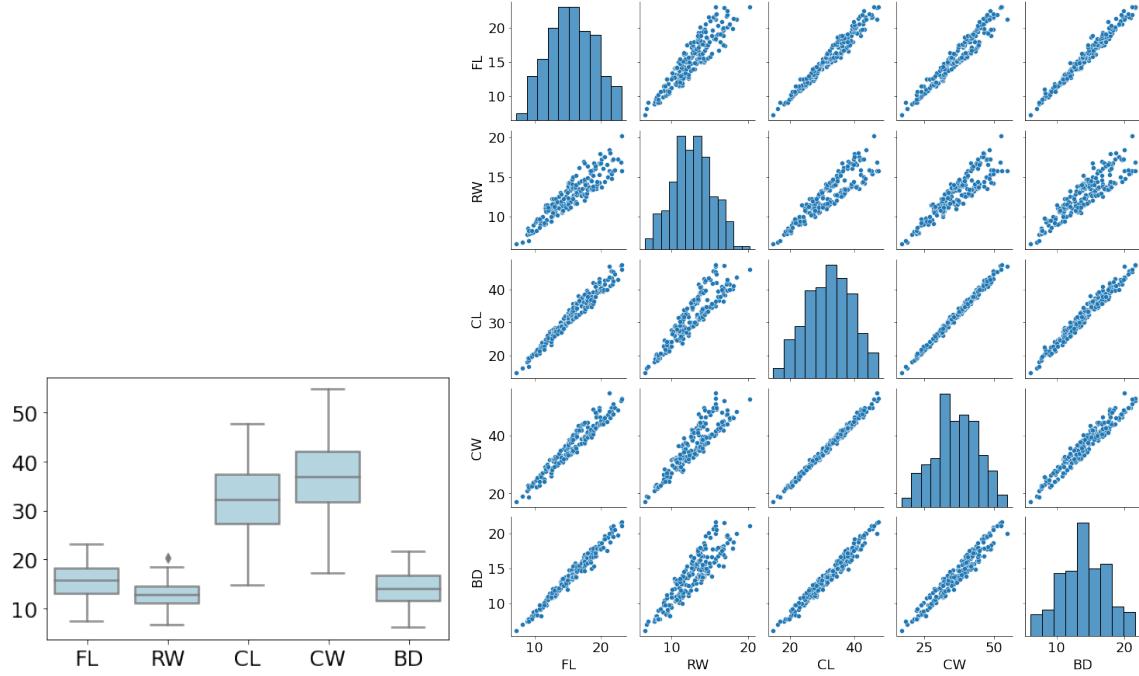


Figure 1.1: (Left) Boxplot and (Right) Pair scatterplot of the crabs dataset

In order to explore this 5-dimensional dataset, we may use exploratory data analysis, and plot some summary statistics. For instance, a boxplot and pair scatterplot are represented in Figure 1.1. The summary plots are useful, but of limited use if the dimensionality  $p$  is high ( $p = 5$  here, but  $p$  is typically a few dozens or even thousands). It is constrained to view data in 2 or 3 dimensions. We may want instead to look for ‘interesting’ projections of  $\mathbf{X}$  into lower dimensions. The hope is that, even though  $p$  is large, considering only carefully selected  $K \ll p$  dimensions is just as informative.

## 1.2 Unsupervised learning as empirical risk minimisation

### 1.2.1 Encoder, decoder and autoencoder

For a data item  $x_i \in \mathbb{R}^p$ , we want to find a lower dimensional representation  $z_i \in \mathbb{R}^K$  with  $K \ll p$ . Such representation should capture important aspects of the data. We denote by  $\text{enc}_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^K$  such transformation, called an **encoder**. The encoder is parameterised by a vector  $\theta \in \Theta$ , and we have

$$\underline{z_i = \text{enc}_\theta(x_i)}.$$

Similarly, we assume that there is another transformation  $\text{dec}_\theta : \mathbb{R}^K \rightarrow \mathbb{R}^p$ , called a **decoder**, that transforms the latent representation  $z_i$  back into the original space  $\mathbb{R}^p$ . This transformation is also parameterised by  $\theta$ . Denote

$$\widehat{x}_i = \underline{\text{dec}_\theta(z_i)}$$

the **reconstructed example**. Let  $h_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^p$  be the function defined by

$$\underline{h_\theta(x_i) = \text{dec}_\theta(\text{enc}_\theta(x_i)) = \widehat{x}_i}$$

and called an **autoencoder**<sup>1</sup>. Note that, as  $K$  is typically smaller than  $p$ , the auto-encoder transformation in general incurs some loss of information and  $\widehat{x}_i \neq x_i$ .

**Example 1.** As a toy example, for  $\theta \in \{1, \dots, p\}$ ,  $K = 1$ , consider the encoder and decoder

$$\begin{aligned} z_i &= \text{enc}_\theta(x_i) = x_{i\theta} \\ \text{dec}_\theta(z_i) &= (0, 0, \dots, \underbrace{z_i}_{\text{coordinate } \theta}, 0, \dots, 0)^\top \end{aligned} \tag{1.1}$$

<sup>1</sup>In the machine learning literature, the term **autoencoder** often refers to the case where both  $\text{enc}_\theta$  and  $\text{dec}_\theta$  are neural networks. We take a broader definition here, where these transformations need not be neural networks.

the encoder simply picks the coordinate  $\theta$  of the vector  $x_i$ ; the decoder uses the value  $z_i$  as the true value for coordinate  $\theta$ , and replace the other values with 0.

Depending on the unsupervised problem, the object of interest may be the transformed vectors  $z_i$ , the parameters  $\theta$ , the reconstructed signal  $\hat{x}_i$  or the decoding function  $\text{dec}_\theta$ .

- Dimensionality reduction and visualisation: the objective is to reduce the dimension of the data and find a lower-dimensional representation of the data example  $x_i$ ; the object of interest is the latent representation  $z_i \in \mathbb{R}^K$  where  $K \ll p$ .
- Data compression: The matrix  $\mathbf{X}$  has  $n \times p$  elements; this requires a lot of memory if  $p$  is large; if the parameter  $\theta$  of the autoencoder is of dimension  $p_\theta$ , storing  $(z_1, \dots, z_n)$  and  $\theta$  instead of  $\mathbf{X}$  only requires storing  $n \times K + p_\theta$  parameters which is typically smaller than  $n \times p$ .
- Clustering: In clustering,  $z_i \in \{1, \dots, K\}$  represents the cluster membership of the data example  $x_i$  to one of  $K$  different groups.
- Restoration: If we assume that the data  $x_i$  are composed of the signal of interest with some added noise/perturbation,  $z_i$  may be interpreted as the signal part, and the reconstructed example  $\hat{x}_i$  as the filtered/restored example, where the noise has been removed. The object of interest is therefore the autoencoder  $h_\theta$  that filters/restores the example  $x_i$

### 1.2.2 Empirical risk minimisation

In order to quantify the error between the original example  $x_i$  and its reconstruction  $\hat{x}_i$ , we introduce a **loss function**  $L : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}_+$  such that  $L(x_i, \hat{x}_i)$  is a measure of the discrepancy between the example and its reconstruction. For unsupervised learning, we will focus on the squared loss function

$$L(x_i, \hat{x}_i) = \|x_i - \hat{x}_i\|^2 = \sum_{j=1}^p (x_{ij} - \hat{x}_{ij})^2.$$

This loss function will enable us to quantify the accuracy of a given autoencoder  $h_\theta$ , via its risk.

**Definition 2** (Risk (unsupervised learning)). Let  $h_\theta$  be an autoencoder, and  $L$  a loss function. The risk of  $h_\theta$  under the loss  $L$  is given by

$$R(h_\theta) = \mathbb{E}[L(X, h_\theta(X))] \quad (\text{population risk})$$

where the expectation is taken with respect to the random variable  $X$ .

The risk corresponds to the expected reconstruction loss/error of an autoencoder  $h_\theta$ . The optimal parameter  $\theta^*$  under the loss  $L$  is the parameter minimising the risk:

$$\theta^* = \arg \min_{\theta \in \Theta} R(h_\theta).$$

The true distribution of  $X$  being unknown, we cannot compute the true risk  $R$ , nor find  $\theta^*$ . We use instead, the **empirical risk**, defined as

$$\widehat{R}_n(h_\theta) = \frac{1}{n} \sum_{i=1}^n L(x_i, h_\theta(x_i)). \quad \begin{matrix} \text{Now we have a sample } x_1, \dots, x_n \\ \text{of } X \end{matrix} \quad (1.2)$$

The empirical risk is the average reconstruction loss over the dataset. The empirical risk minimiser is defined as

$$\widehat{\theta} = \arg \min_{\theta \in \Theta} \widehat{R}_n(h_\theta).$$

**Example 3** (continued). Consider the same autoencoder as above. The loss is

$$L(x_i, h_\theta(x_i)) = \sum_{j \neq \theta} x_{ij}^2$$

Assume that  $X$  has mean  $\mu$  and covariance matrix  $\Sigma$ . The risk is given by

$$R(h_\theta) = \sum_{j \neq \theta} \mathbb{E}(X_j^2) = \sum_{j \neq \theta} \Sigma_{jj} + \mu_j^2 = \left( \sum_{j=1}^p (\Sigma_{jj} + \mu_j^2) \right) - (\Sigma_{\theta\theta} + \mu_\theta^2)$$

The optimal parameter is

$$\theta^* = \arg \max_{j=1, \dots, p} \Sigma_{jj} + \mu_j^2. \quad \begin{matrix} \text{maximize this} \\ \theta = 1, \dots, p \end{matrix}$$

The empirical risk is

$$\widehat{R}_n(h_\theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j \neq \theta} x_{ij}^2 = \left( \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p x_{ij}^2 \right) - \frac{1}{n} \sum_{i=1}^n x_{i\theta}^2 \quad (1.3)$$

and the empirical risk minimiser is defined as

$$\widehat{\theta} = \arg \max_{j=1, \dots, p} \underbrace{\sum_{i=1}^n x_{ij}^2}_{}$$

As we will see in the next two chapters, both Principal Component Analysis and K-means clustering can be cast under this general framework.

# 2 — Principal Component Analysis

Principal component analysis (PCA) is one of the most widely used tools in data analysis. It is a linear dimensionality reduction technique that aims at finding a new basis to represent a noisy dataset. It is often used as a preprocessing tool to reduce the dimension before supervised techniques are used.

For simplicity, we will assume that the data are centered, that is  $\sum_{i=1}^n x_{ij} = 0$  for all  $j$ .

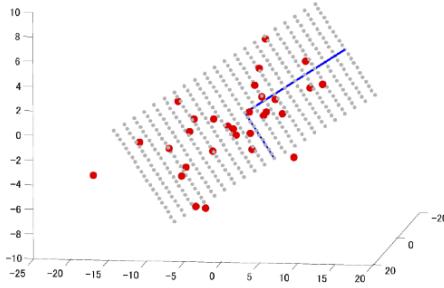


Figure 2.1: Representation of the first two scaled principal components (blue) on a three-dimensional dataset (red).

## 2.1 Overview

We first give an overview of PCA and of its two interpretations as variance maximisation and empirical risk minimisation. PCA finds an orthogonal basis (principal components)  $v_1, v_2, \dots, v_p$  such that:

- The first principal component (PC)  $v_1$  is the *direction of greatest variance* of the data.
- The  $j$ -th PC  $v_j$  is the *direction orthogonal to  $v_1, v_2, \dots, v_{j-1}$  of greatest variance*, for  $j = 2, \dots, p$ .

The  $K$ -dimensional representation of the data item  $x_i$  is the vector of projections of  $x_i$  onto the first  $K$  PCs:

$$z_i = \underbrace{V_{1:K}^\top x_i}_{\text{where } V_{1:K} = (v_1, \dots, v_K)} = (v_1^\top x_i, \dots, v_K^\top x_i)^\top \in \mathbb{R}^K, \quad \left( \begin{array}{c} v_1^\top \\ \vdots \\ v_K^\top \end{array} \right) \left( \begin{array}{c} x_i \\ \vdots \\ x_i \end{array} \right)$$

For a general orthonormal  $p$ -by- $K$  matrix  $A$ , denote  $\text{enc}_A$  and  $\text{dec}_A$  the linear encoder and decoder defined by

$$\begin{aligned} z_i &= \text{enc}_A(x_i) = A^\top x_i \\ x_i &= \text{dec}_A(z_i) = Az_i. \end{aligned}$$

Additionally, let  $h_A(x_i) = AA^\top x_i$  be the associated linear autoencoder. The matrix  $V_{1:K}$  is an empirical risk minimiser for the class of linear autoencoder  $h_A$  under the squared loss:

$$V_{1:K} \in \arg \min_{A \in \mathcal{A}} \frac{1}{n} \sum_{i=1}^n \|x_i - h_A(x_i)\|^2 \quad (2.1)$$

where  $\mathcal{A}$  is the set of  $p$ -by- $K$  orthonormal matrices.

## 2.2 Eigenvalue decomposition of the covariance matrix and sample covariance matrix.

Our dataset is an i.i.d. sample  $(x_1, \dots, x_n)$  of a random vector  $X = (X_1 \dots X_p)^\top$ . Let  $\Sigma$  denote the  $p$ -by- $p$  covariance matrix of the random vector  $X$ .

**Covariance matrix.** The covariance matrix  $\Sigma$  is a

- **real and symmetric** matrix, so there exist  $p$  eigenvectors  $v_1^*, \dots, v_p^*$  that are pairwise orthogonal and  $p$  associated eigenvalues  $\lambda_1^* \geq \lambda_2^* \dots \geq \lambda_p^*$  which satisfy the eigenvalue equation  $\Sigma v_i^* = \lambda_i^* v_i^*$ . In particular, the  $p$ -by- $p$  matrix  $V^* = (v_1^*, \dots, v_p^*)$  is an orthogonal matrix:

$$\underbrace{V^*(V^*)^\top}_{= (V^*)^\top V^*} = I_p.$$

- **positive-semidefinite** matrix, so the eigenvalues are non-negative:

$$\underbrace{\lambda_i^* \geq 0, \forall i.}$$

The **eigenvalue decomposition** of the covariance matrix is given by

$$\underbrace{\Sigma = V^* \Lambda^* (V^*)^\top}$$

where  $\Lambda^*$  is a diagonal matrix with eigenvalues

$$\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_p^* \geq 0.$$

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

The total variance of the random vector  $X$  is

$$\begin{aligned} \underbrace{\text{TV}(X)}_{=} &= \mathbb{E} \left[ \sum_{j=1}^p (X_j - \mathbb{E}[X_j])^2 \right] = \sum_{j=1}^p \Sigma_{jj} = \underbrace{\text{trace}(\Sigma)}_{=} \\ &= \underbrace{\text{trace}(V^* \Lambda^* (V^*)^\top)}_{=} = \text{trace}((V^*)^\top V^* \Lambda^*) = \underbrace{\text{trace}(\Lambda^*)}_{=} = \sum_{j=1}^p \lambda_j^* \end{aligned}$$

**Sample covariance matrix.** Similarly, let  $S$  be the sample covariance matrix of the data. Assuming the data have been centered, that is  $\sum_{i=1}^n x_{ij} = 0$  for all  $j$ , it is defined as

$$\underbrace{S}_{=} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^\top = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}.$$

The sample covariance matrix  $S$  is a

- **real and symmetric** matrix, so there exist  $p$  eigenvectors  $v_1, \dots, v_p$  that are pairwise orthogonal and  $p$  associated eigenvalues  $\lambda_1, \dots, \lambda_p$  which satisfy the eigenvalue equation  $S v_i = \lambda_i v_i$ . In particular,  $V$  is an orthogonal matrix:

$$\underbrace{V V^\top}_{=} = V^\top V = I_p.$$

- **positive-semidefinite** matrix, so the eigenvalues are non-negative:

$$\lambda_i \geq 0, \forall i.$$

The **eigenvalue decomposition** of the sample covariance matrix  $S$  is given by

$$\underbrace{S = V \Lambda V^\top}_{=}$$

where  $\Lambda$  is a diagonal matrix with eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0.$$

The sample total variance of the data  $\mathbf{X}$  is

$$\begin{aligned} \underbrace{\text{TV}(\mathbf{X})}_{=} &= \frac{1}{n-1} \sum_{i=1}^n \sum_{j=1}^p x_{ij}^2 = \text{trace}(S) \\ &= \text{trace}(V \Lambda V^\top) = \text{trace}(\Lambda) = \sum_{j=1}^p \lambda_j \end{aligned}$$

## 2.3 Derivation of the PCs

We first derive the population PCs, that is the directions of maximum variance of the random vector  $X$ . We then proceed with the PCs of the data  $\mathbf{X}$ .

### 2.3.1 Population PCs

Assume for simplicity that  $X$  has zero mean.

$$(a_1 \quad a_p) \begin{pmatrix} \vdots & x_1^\top \\ \vdots & \vdots \\ \vdots & x_p^\top \end{pmatrix}$$

**First PC.** For the 1<sup>st</sup> PC, we seek a derived scalar variable  $Z_1$  of the form

$$Z_1 = a_1^\top X = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

where  $a_1 = (a_{11}, \dots, a_{1p})^\top \in \mathbb{R}^p$  is a unit-norm vector, that is  $\underline{a_1^\top a_1 = 1}$ . For any fixed vector  $a_1$ ,

$$\var(Z_1) = \var(a_1^\top X) = a_1^\top \Sigma a_1.$$

where  $\Sigma$  is the covariance matrix of the random vector  $X$ . The (population) first principal component  $v_1^*$  is the unit-norm vector  $a_1$  that maximises  $\var(Z_1)$ , that is

$$\underbrace{v_1^* = \arg \max_{a_1 \in \mathbb{R}^p} a_1^\top \Sigma a_1}_{\text{subject to: } \underline{a_1^\top a_1 = 1}}.$$

The Lagrangian of the problem is given by:

$$\mathcal{L}(a_1, \gamma_1) = \underline{a_1^\top \Sigma a_1 - \gamma_1(a_1^\top a_1 - 1)}.$$

The corresponding vector of partial derivatives is

$$\frac{\partial \mathcal{L}(a_1, \gamma_1)}{\partial a_1} = 2\Sigma a_1 - 2\gamma_1 a_1.$$

Setting this to zero gives  $\underline{\Sigma a_1 = \gamma_1 a_1}$ . We recognize the eigenvector equation, i.e.  $a_1$  must be an eigenvector of  $\Sigma$  and the dual variable  $\gamma_1$  is the corresponding eigenvalue. Since  $a_1^\top \Sigma a_1 = \gamma_1 a_1^\top a_1 = \underline{\gamma_1}$ , the first PC must be the eigenvector  $v_1^*$  associated with the largest eigenvalue  $\lambda_1^*$  of  $\Sigma$ . We therefore have

$$\var(Z_1) = \underline{\lambda_1^*}.$$

**Second and subsequent PCs.** To derive the subsequent principal components, we proceed as before. We consider the scalar variable

$$Z_2 = a_2^\top X$$

where  $a_2 \in \mathbb{R}^p$  is a unit-norm vector, that is  $\underline{a_2^\top a_2 = 1}$ , which is orthogonal to the first PC:  $a_2^\top v_1^* = 0$ . We wish to find the vector  $a_2$  that maximises  $\var(Z_2) = a_2^\top \Sigma a_2$ . Let

$$\underbrace{v_2^* = \arg \max_{a_2 \in \mathbb{R}^p} a_2^\top \Sigma a_2}_{\text{subject to: } \underline{a_2^\top a_2 = 1} \text{ and } \underline{a_2^\top v_1^* = 0}}.$$

The Lagrangian of the problem is given by

$$\mathcal{L}(a_2, \gamma_2, \mu) = \underline{a_2^\top \Sigma a_2 - \gamma_2(a_2^\top a_2 - 1) - \mu(v_1^*)^\top a_2}.$$

The corresponding vector of partial derivatives is

$$\frac{\partial \mathcal{L}(a_2, \gamma_2, \mu)}{\partial a_2} = \underline{2\Sigma a_2 - 2\gamma_2 a_2 - \mu v_1^*} = 0.$$

Left-multiplying the above expression by  $a_2^\top$  gives

$$a_2^\top \Sigma a_2 = \underline{\gamma_2}.$$

Multiplying by  $a_2$  gives

$$\Sigma a_2 = \underline{\gamma_2 a_2}.$$

Hence  $a_2$  must be an eigenvector of  $\Sigma$  with corresponding eigenvalue  $\gamma_2$ , where  $\gamma_2$  is the quantity we wish to maximise. The second PC is therefore the second eigenvector  $v_2^*$  of  $\Sigma$ , with corresponding second highest eigenvalue  $\lambda_2^*$ . We have

$$\var(Z_2) = \underline{\lambda_2^*}.$$

By induction, we conclude that the population PCs  $v_1^*, \dots, v_p^*$  are the  $p$  eigenvectors of  $\Sigma$ , with associated eigenvalues  $\lambda_1^* \geq \lambda_2^* \geq \dots \geq \lambda_p^* \geq 0$ . The PCs are orthogonal.

### 2.3.2 Empirical PCs.

In practice, we do not know the *true* covariance matrix  $\Sigma$ , so cannot compute the population PC. We want instead to find projections maximising the sample variance. For a vector  $a \in \mathbb{R}^p$ , the projections  $(a^\top x_1, \dots, a^\top x_n)$  of the examples  $x_i$  onto  $a$  have sample variance  $a^\top S a$ . We therefore proceed as for the population PCs, replacing  $\Sigma$  with the sample covariance matrix  $S$ . The first two empirical PCs  $v_1$  and  $v_2$  are therefore given by

$$\begin{aligned} v_1 &= \arg \max_{a_1 \in \mathbb{R}^p} a_1^\top S a_1 \quad \text{subject to: } a_1^\top a_1 = 1 \\ v_2 &= \arg \max_{a_2 \in \mathbb{R}^p} a_2^\top S a_2 \quad \text{subject to: } a_2^\top a_2 = 1 \text{ and } v_1^\top a_2 = 0. \end{aligned}$$

and so on, where the PCs  $v_1, v_2, \dots, v_p$  are the eigenvectors of the sample covariance matrix  $S$ , with corresponding eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_p$ . The projection of the  $i$ th example  $x_i$  onto the  $j$ th PC component is therefore

$$z_{ij} = \underline{v_j^\top x_i}$$

and we note  $z_i = (z_{i1}, \dots, z_{iK})^\top$  the  $K$  dimensional latent representation of the example  $x_i$ .

## 2.4 Properties

The projections onto the principal components have variance/sample variance given by the eigenvalues of  $\Sigma/S$ :

$$\begin{aligned} \text{var}(Z_j) &= \text{var}((v_j^*)^\top X) = \underline{\lambda_j^*} \\ &\underbrace{\frac{1}{n-1} \sum_{i=1}^n z_{ij}^2}_{= \lambda_j}. \end{aligned}$$

The projections onto the principal components are uncorrelated: for  $j \neq k$ ,

$$\begin{aligned} \text{cov}(Z_j, Z_k) &= (v_j^*)^\top \Sigma v_k^* = \lambda_k^* (v_j^*)^\top v_k^* = 0 \\ &\frac{1}{n-1} \sum_{i=1}^n z_{ij} z_{ik} = v_j^\top S v_k = \lambda_k v_j^\top v_k = 0. \end{aligned}$$

The total variance of  $X$  is  $\text{TV}(X) = \sum_{i=1}^p \Sigma_{ii} = \sum_{j=1}^p \lambda_j^*$ . The total variance explained by the  $j^{th}$  PC is  $\text{TV}(Z_j) = \lambda_j^*$ . The proportion of total variance explained by the first  $K$  (population) PCs is therefore

$$\frac{\sum_{j=1}^K \lambda_j^*}{\sum_{\ell=1}^p \lambda_\ell^*}.$$

Similarly, the sample total variance is  $\sum_i S_{ii} = \overbrace{\sum_{j=1}^p \lambda_j}$ . The proportion of total variance explained by the first  $K$  PCs is

$$\frac{\sum_{j=1}^K \lambda_j}{\sum_{\ell=1}^p \lambda_\ell}.$$

## 2.5 PCA as empirical risk minimisation

Let  $h_A(x_i) = AA^\top x_i$  be the autoencoder parameterised by an orthonormal  $p$ -by- $K$  matrix  $A$  (that is  $A^\top A = I_K$ ). The risk of this autoencoder for a squared loss function is

$$R(h_A) = \mathbb{E}[||X - AA^\top X||^2]$$

where the expectation is taken with respect to the random vector  $X$ . The empirical risk is

$$\widehat{R}_n(h_A) = \frac{1}{n} \sum_{i=1}^n ||x_i - AA^\top x_i||^2.$$

  
**Proposition 4.** Let  $V_{1:K}^* = (v_1^*, \dots, v_K^*)$  be  $p$ -by- $K$  matrix of the first  $K$  population PCs. Then  $V_{1:K}^*$  minimises the risk  $R(h_A)$  amongst the set of  $p$ -by- $K$  orthonormal matrices  $A$ .

**Proposition 5.** Let  $V_{1:K} = (v_1, \dots, v_K)$  be  $p$ -by- $K$  matrix of the first  $K$  PCs. Then  $V_{1:K}$  minimises the empirical risk  $\widehat{R}_n(h_A)$  amongst the set of  $p$ -by- $K$  orthonormal matrices  $A$ .

Proof: See problem sheet.

**Example 6** (Crabs). We apply PCA to the crabs dataset, with  $K = 5$ . The estimated principal components and proportion of variance explained are as follows:

Principal Components/Loadings:

$$\text{FL} \begin{pmatrix} -0.289 & -0.323 & 0.507 & 0.734 & 0.125 \\ -0.197 & -0.865 & -0.414 & -0.148 & -0.141 \\ -0.599 & 0.198 & 0.175 & -0.144 & -0.742 \\ -0.662 & 0.288 & -0.491 & 0.126 & 0.471 \\ -0.284 & -0.160 & 0.547 & -0.634 & 0.439 \end{pmatrix} = V$$

$v_1 \quad v_2 \quad v_3 \quad v_4 \quad v_5$

Importance of components:

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5
Proportion of Variance	0.9824718	0.009055108	0.006984337	0.0009447218	0.0005440328
Cumulative Proportion	0.9824718	0.991526908	0.998511245	0.9994559672	1.00000000000

The first PC has the same sign for all variables, and it is related to the size of the crabs. It explains most of the variance in the data. Is this latent representation useful? The dataset also contains some class information with 4 classes (species B/O and sex M/F): BF, BM, OF, OM. Plotting a pair scatterplot of the transformed features  $z_i$  reveals some useful structure: when looking at the second and third PC coordinates, we obtain a good separation of the classes. Applying some classifier on the transformed features  $(z_{i2}, z_{i3})$  is likely to provide some good performance. Note that using the first two components does not provide such a good separation for the classes, see Figure 6

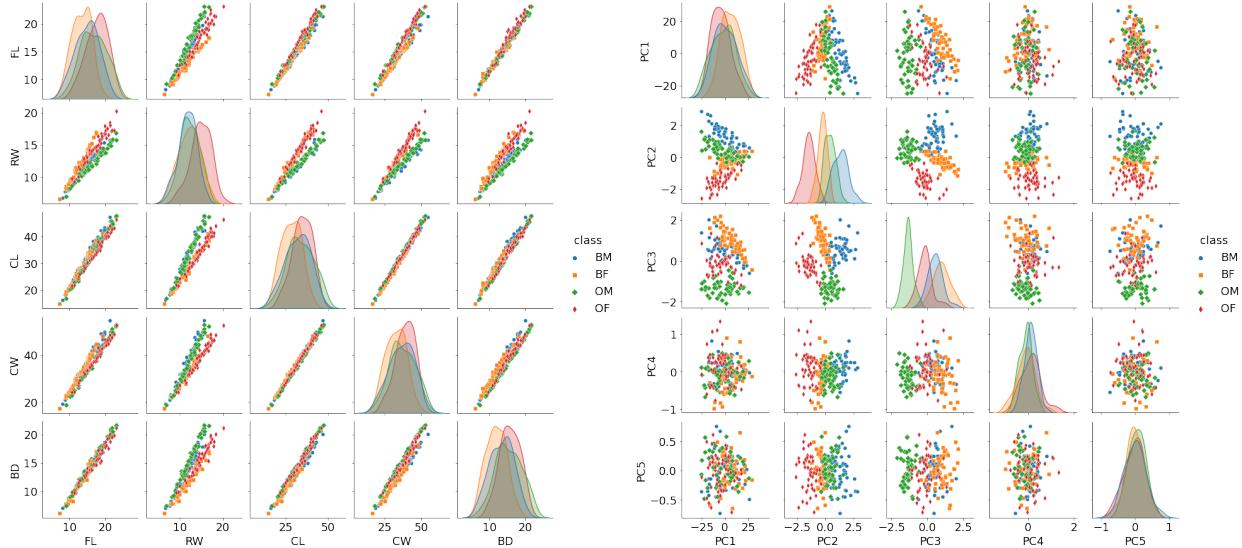


Figure 2.2: Pair scatterplot on the crabs dataset on the original data  $x_i$  (left) and on the transformed PCA features  $z_i$  (right), with class information.

**Example 7** (Eigenfaces). We consider an application of PCA to dimension reduction and compression for image data. The Olivetti dataset contains 400 64-by-64 images, quantized to 256 grey levels. The dimension is therefore  $p = 64 \times 64 = 4096$ , and  $n = 400$ . A sample of the observations is given in Figure 7.

In this application, the estimated principal components/eigenvectors  $v_1, v_2, \dots$ , are called eigenfaces. They are used as a dimensionality reduction technique. The eigenface projections  $z_i$  can be used as low-dimensional features in pattern recognition algorithms for image classification. When using a small set  $K$  of eigenfaces, each image can be reconstructed (with some loss) using a weighted sum of the eigenfaces, thus leading to less parameters to store ( $K(n + p)$  versus  $np$ ). The first estimated eigenfaces, and the reconstructed images for a subset of the images are shown in Figure 7.

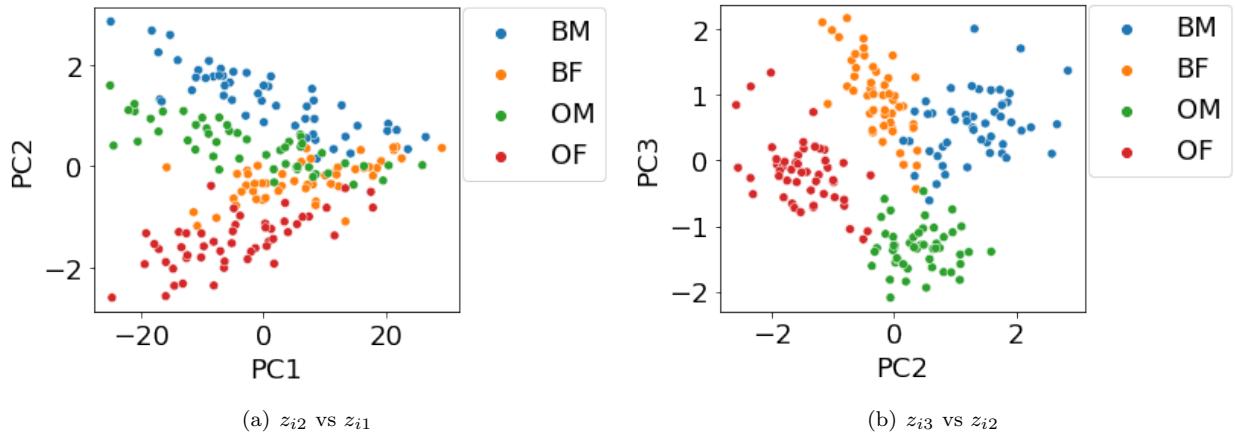


Figure 2.3: Scatterplot of the transformed features, with class information.



Figure 2.4: Sample images from the Olivetti faces dataset.

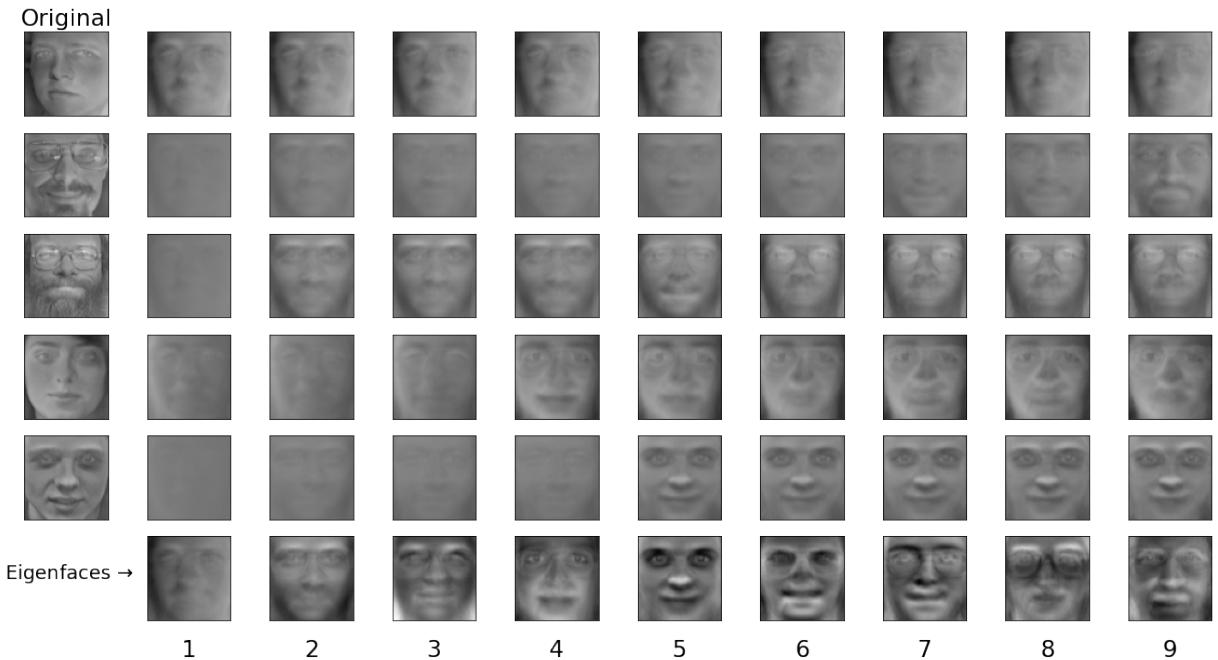


Figure 2.5: Olivetti faces dataset. Bottom row: First 9 eigenfaces. Rows 1-5: Original image, followed by the reconstructed image using 1 to 9 eigenfaces.

## 2.6 Comments on the use of PCA

PCA is commonly used to project data  $\mathbf{X}$  onto the first  $K$  PCs giving the  $K$ -dimensional view of the data that best preserves the first two moments. Although PCs are uncorrelated, scatterplots sometimes reveal structures

in the data other than linear correlation. The emphasis on the variance is where the weaknesses of PCA stem from:

- Assuming large variances are meaningful (high signal-to-noise ratio)
- The PCs depend heavily on the units measurement. Where the data matrix contains measurements of vastly differing orders of magnitude, the PC will be greatly biased in the direction of larger measurement. In these cases, it is recommended to calculate PCs from  $\text{Corr}(X)$  instead of  $\text{cov}(X)$
- Lack of robustness to outliers: the variance is affected by outliers and so are PCs.
- Sample size (e.g. the number of crabs collected for each sub-species) will have an effect on the PCs.

## 2.7 SVD and PCA

Any real-valued  $n \times p$  matrix  $\mathbf{X}$  can be written as  $X = UDV^\top$  where

- $U$  is an  $n \times n$  orthogonal matrix:  $UU^\top = U^\top U = I_n$
- $D$  is a  $n \times p$  matrix with decreasing non-negative elements on the diagonal (the singular values) and zero off-diagonal elements.
- $V$  is a  $p \times p$  orthogonal matrix:  $VV^\top = V^\top V = I_p$

This is known as the **Singular Value Decomposition (SVD)** of  $\mathbf{X}$ . The SVD always exists, even for non-square matrices. Fast and numerically stable algorithms for SVD are available in most packages. Let  $\mathbf{X} = UDV^\top$  be the SVD of the  $n \times p$  data matrix  $\mathbf{X}$ . Note that

$$(n-1)S = \underbrace{\mathbf{X}^\top \mathbf{X}}_{p \times p} = (UDV^\top)^\top (UDV^\top) = VD^\top U^\top UDV^\top = VD^\top DV^\top,$$

using orthogonality ( $U^\top U = I_n$ ) of  $U$ . The eigenvalues of  $S$  are thus the diagonal entries of  $\Lambda = \frac{1}{n-1} D^\top D$ . We also have (using orthogonality  $V^\top V = I_p$ )

$$\underbrace{\mathbf{X} \mathbf{X}^\top}_{n \times n} = (UDV^\top)(UDV^\top)^\top = UDV^\top VD^\top U^\top = \underbrace{UDD^\top U^\top}_{n \times n},$$

The  $n$ -by- $n$  matrix

$$\mathbf{B} = \underbrace{\mathbf{X} \mathbf{X}^\top}_{n \times n},$$

such that  $\mathbf{B}_{ij} = x_i^\top x_j$ , is called the **Gram matrix** of the dataset  $\mathbf{X}$ .  $\mathbf{B}$  and  $(n-1)S = \mathbf{X}^\top \mathbf{X}$  have the same nonzero eigenvalues, equal to the non-zero squared singular values of  $\mathbf{X}$ . The transformed variables

$$\mathbf{Z} = \underbrace{\mathbf{X} V}_{n \times p} = UDV^\top V = UD$$

can be obtained by eigendecomposition of  $\mathbf{B}$ , which requires less computation than calculating the eigendecomposition of  $S$  if  $p > n$ .

$$\mathbf{D}\mathbf{D}^\top : n \times n.$$

$$\begin{aligned} S &= \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} \\ &= \frac{1}{n-1} \mathbf{V} \mathbf{D}^\top \mathbf{D} \mathbf{V}^\top \\ &= \mathbf{V} \Delta \mathbf{V}^\top \end{aligned}$$

$$\Delta = \frac{\mathbf{D}^\top \mathbf{D}}{n-1}$$

# 3 — K-means clustering

## 3.1 Clustering

Many datasets consist of multiple heterogeneous subsets. The aim of *cluster analysis* is, given unlabelled data, to automatically group the data examples into coherent subsets/clusters. Examples include:

- market segmentation of shoppers based on browsing and purchase histories
- discovery of different types of breast cancer based on the gene expression measurements
- discovery of communities in social networks
- image segmentation

**Definition 8. (Partition)** Let  $n \geq 1$  an integer. A **partition**  $\Pi = \{C_1, \dots, C_K\}$  of the set of integers  $\{1, \dots, n\}$  is such that

- For all  $k = 1, \dots, K$ ,  $C_k \neq \emptyset$  and  $C_k \subseteq \{1, \dots, n\}$
- For all  $k \neq k'$ ,  $C_k \cap C_{k'} = \emptyset$
- $\bigcup_{k=1}^K C_k = \{1, \dots, n\}$ .

$C_k$  is known as the  $k$ th group or cluster, of size  $1 \leq |C_k| \leq n$ .  $K$  is the number of groups/clusters.

For instance

$$\Pi = \underbrace{\{\{1, 2, 5\}\}}_{C_1}, \underbrace{\{\{3, 4, 7\}\}}_{C_2}, \underbrace{\{\{6\}\}}_{C_3}$$

is a partition of  $\{1, 2, 3, 4, 5, 6, 7\}$  with  $K = 3$  groups. A partition can be encoded by **cluster labels**  $(z_1, \dots, z_n)$ , where  $z_i \in \{1, \dots, K\}$  is such that  $\underline{z_i = k}$  if  $i \in C_k$ . For the example above, the cluster labels are

$$z_1 = z_2 = z_5 = 1, \quad z_3 = z_4 = z_7 = 2, \quad z_6 = 3.$$

There are two broad classes of clustering methods: model-free clustering and model-based clustering. In model-based clustering, each cluster is described using a probability model. We will not look at this type of approach in this course. Model-free clustering relies on the notion of dissimilarity between pairs of data examples in a cluster.

**Model-free clustering.** Let  $d = (x_1, \dots, x_n)$  be a dataset of size  $n$ , where  $x_i \in \mathbb{R}^p$  and  $\rho : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}_+$  be a **dissimilarity function**. For two data examples  $x_i$  and  $x_j$ ,  $\rho(x_i, x_j)$  is a measure of how dissimilar the two examples are. A typical dissimilarity function is the squared Euclidian distance

$$\rho(x_i, x_j) = \|x_i - x_j\|^2.$$

The notion of dissimilarity between data items is central. Its choice will depend on the dataset being analysed and dictated by domain specific knowledge. A model-free clustering method is a map  $\mathcal{F} : (d, \rho) \mapsto \{C_1, \dots, C_K\}$  which takes as input a dataset  $d$  of size  $n$  and a dissimilarity function  $\rho$  and returns a partition of  $\{1, \dots, n\}$ . Intuitively, clustering aims to group similar items together and to place dissimilar items into different groups. The two objectives can contradict each other (similarity is not a transitive relation, while being in the same cluster is an equivalence relation). Three basic properties are desired:

- **Scale invariance.** For any  $\alpha > 0$ ,  $\mathcal{F}(d, \alpha\rho) = \mathcal{F}(d, \rho)$ .
- **Richness.** For any partition  $\Pi = \{C_1, \dots, C_K\}$  of  $\{1, \dots, n\}$ , there exists dissimilarity  $\rho$ , such that  $\mathcal{F}(d, \rho) = \Pi$ .
- **Consistency.** If  $\rho$  and  $\rho'$  are two dissimilarities such that for all  $x_i, x_j$  the following holds:

$$\begin{aligned} x_i, x_j \text{ belong to the same cluster in } \mathcal{F}(d, \rho) &\implies \rho'(x_i, x_j) \leq \rho(x_i, x_j) \\ x_i, x_j \text{ belong to different clusters in } \mathcal{F}(d, \rho) &\implies \rho'(x_i, x_j) \geq \rho(x_i, x_j), \end{aligned}$$

then  $\mathcal{F}(d, \rho') = \mathcal{F}(d, \rho)$ .

Kleinberg (2003) proved that there exists no clustering method that satisfies all three properties. Every algorithm therefore has to find a trade-off between the above three properties.

## 3.2 K-means

K-means is a model-free clustering method whose goal is to partition  $n$  data items  $x_1, \dots, x_n$ , with  $x_i \in \mathbb{R}^p$ , into a number  $K$  of clusters  $C_1, \dots, C_K$ , where  $K$  is pre-defined. We use as a dissimilarity function the squared Euclidian distance

$$\rho(x_i, x_{i'}) = \|x_i - x_{i'}\|^2 = \sum_{j=1}^p (x_{ij} - x_{i'j})^2.$$

The dissimilarity of a cluster  $C_k$  will be expressed as the sum of the dissimilarities of pairs of data in that cluster. We consider the following overall objective function, obtained by summing the dissimilarities of each cluster  $C_k$ , and scaled by  $1/2n$

$$\widetilde{W}(C_1, \dots, C_K) = \frac{1}{2n} \sum_{k=1}^K \sum_{i, i' \in C_k} \rho(x_i, x_{i'}) = \frac{1}{2n} \sum_{k=1}^K \sum_{i, i' \in C_k} \|x_i - x_{i'}\|^2.$$

$\widetilde{W}(C_1, \dots, C_K)$  will be called the **K-means objective function**. We want to find the partition  $\{C_1, \dots, C_K\}$  that minimises  $\widetilde{W}$ . The total number of partitions of  $n$  elements in  $K$  clusters is  $S(n, K)$ , the Stirling number of the second kind. This number increases very quickly with  $n$ . For example, for  $n = 100$  and  $K = 4$ , we have

$$S(100, 4) = 66955751844038698560793085292692610900187911879206859351901$$

different partitions. Exhaustive search to minimise  $\widetilde{W}$  is therefore practically impossible. The K-means algorithm is an iterative algorithm that allows to find a local minimum of the objective function  $\widetilde{W}$ .

### 3.2.1 Extended K-means objective function

We first show that the partition minimising  $\widetilde{W}$  is also the minimiser of another objective function  $W$  defined on an extended space. Let

$$\bar{x}_{C_k} = \underbrace{\frac{1}{|C_k|} \sum_{i \in C_k} x_i}_{\text{be the empirical mean of the data in cluster } C_k}. \quad \text{HS: } \frac{1}{2n} \sum_{i, i'} \|x_i - x_{i'}\|^2 = \frac{1}{2n} \sum_{i, i'} \|x_i - \bar{x}_{C_k} + \bar{x}_{C_k} - x_{i'}\|^2$$

Additionally, note that

$$\bar{x}_{C_k} = \arg \min_{\mu_k \in \mathbb{R}^p} \sum_{i \in C_k} \|x_i - \mu_k\|^2.$$

Hence we can write

$$\widetilde{W}(C_1, \dots, C_K) = \min_{\mu_1, \dots, \mu_K} \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2.$$

Consider now the joint optimisation problem over the partition  $\{C_1, \dots, C_K\}$  and **prototypes** or **cluster centroid**  $(\mu_1, \dots, \mu_K)$ , with extended objective function

$$W(C_1, \dots, C_K, \mu_1, \dots, \mu_K) = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2.$$

Define

$$(\hat{C}_1, \dots, \hat{C}_K, \hat{\mu}_1, \dots, \hat{\mu}_K) = \arg \min_{C_1, \dots, C_K, \mu_1, \dots, \mu_K} W(C_1, \dots, C_K, \mu_1, \dots, \mu_K).$$

The value  $(\hat{C}_1, \dots, \hat{C}_K)$  is therefore also a minimiser of the K-means objective  $\widetilde{W}$ .

Joint minimisation of  $W$  over both the partition and the cluster centroids is still computationally difficult. However, minimising  $W$  with respect the partition when the centroids are fixed, or minimising  $W$  with respect to the centroids when the partition is fixed are both easy tasks as described below. Note that

$$W = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2 = \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

where  $z_i = k$  if and only if  $i \in C_k$ .

$$W = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2 = \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

- Given the partition  $\{C_1, \dots, C_K\}$ , we can find the optimal prototypes easily by differentiating  $W$  with respect to  $\mu_k$ :

$$\frac{\partial W}{\partial \mu_k} = -2 \sum_{i \in C_k} (x_i - \mu_k) = 0 \Rightarrow \mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

- Given the prototypes  $(\mu_1, \dots, \mu_K)$ , we can easily find the optimal partition/cluster labels by assigning each data point to the closest cluster prototype:

$$z_i = \arg \min_{k=1, \dots, K} \|x_i - \mu_k\|^2$$

This suggests the use a coordinate descent optimisation algorithm, which optimises with respect to one coordinate at a time, while keeping the other coordinate fixed.

### Algorithm 1 Coordinate descent algorithm

Let  $f(\theta_1, \theta_2) \in \mathbb{R}$  be some objective function we want to minimise with respect to  $(\theta_1, \theta_2)$ . The coordinate descent algorithm proceeds as follows:

- Initialise  $\theta_2^{(0)}$  to some value and set  $t = 0$
- Repeat until convergence
  - Set  $t \leftarrow t + 1$
  - $\theta_1^{(t)} = \arg \min_{\theta_1} f(\theta_1, \theta_2^{(t-1)})$
  - $\theta_2^{(t)} = \arg \min_{\theta_2} f(\theta_1^{(t)}, \theta_2)$
- Return  $(\theta_1^{(t)}, \theta_2^{(t)})$

Each iteration of the algorithm decreases the value of the objective function:  $f(\theta_1^{(1)}, \theta_2^{(1)}) \geq f(\theta_1^{(2)}, \theta_2^{(2)}) \geq f(\theta_1^{(3)}, \theta_2^{(3)}) \dots$

### 3.2.2 K-means algorithm

The K-means algorithm is a widely used method that returns a *local optimum* of the objective function  $W$ , using coordinate descent.

### Algorithm 2 K-means algorithm

- Randomly initialise  $K$  cluster centroids  $\mu_1, \dots, \mu_K$ .
- Cluster assignment: For each  $i = 1, \dots, n$ , assign each  $x_i$  to the cluster with the nearest centroid,

$$z_i := \arg \min_k \|x_i - \mu_k\|^2$$

Set  $C_k := \{i : z_i = k\}$  for each  $k$ .

- Move centroids: Set  $\mu_1, \dots, \mu_K$  to the averages of the new clusters:

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

- Repeat steps 2-3 until convergence.
- Return the partition  $\{C_1, \dots, C_K\}$  and means  $\mu_1, \dots, \mu_K$ .

The algorithm stops in a finite number of iterations. Between steps 2 and 3,  $W$  either stays constant (same partitions) or it decreases, this implies that we never revisit the same partition. As there are only finitely many partitions, the number of iterations cannot exceed this. The K-means algorithm need not converge to a global optimum. K-means can get stuck at suboptimal configurations and the result depends on the starting configuration. Typically, we perform a number of runs with different initial configurations, and pick the end result with minimum  $W$ . The result of K-means algorithm with different initialisations is shown in Figure 3.2

### 3.2.3 K-means as empirical risk minimisation

The minimiser of the extended K-means objective can also be interpreted in terms of empirical risk minimisation. Recall that  $(x_1, \dots, x_n)$  are iid realisations of a random variable  $X$ . For a parameter  $\theta = (\mu_1, \dots, \mu_K)$ , define

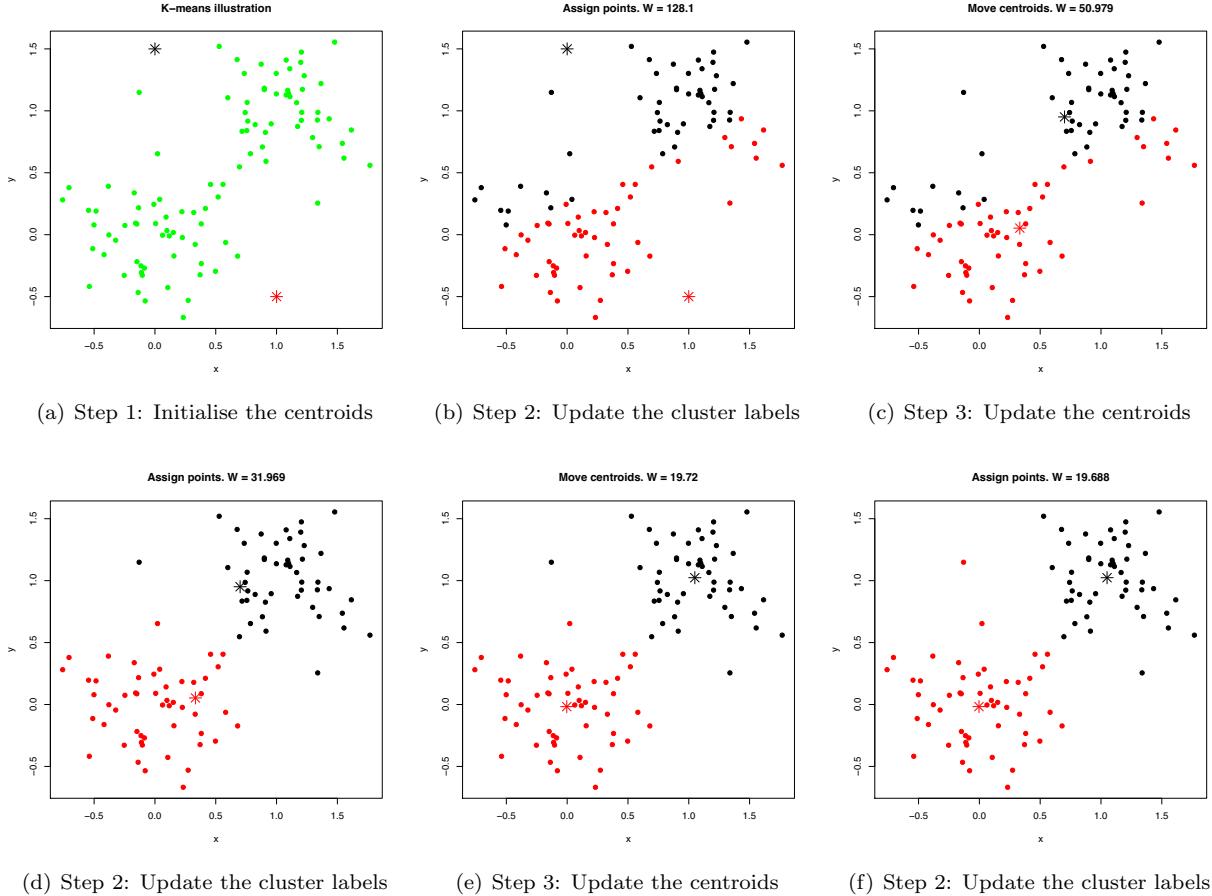


Figure 3.1: Illustration of the first iterations of the K-means algorithms for  $x_i \in \mathbb{R}^2$ . Observations are represented by dots. The two cluster means are represented by a red and a black stars. The cluster assignments of the data examples are indicated by their color (red or black).

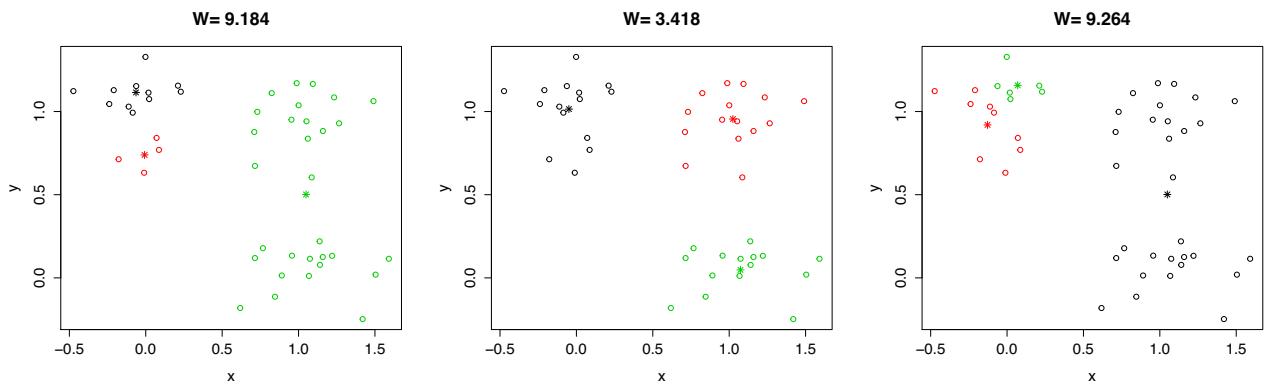


Figure 3.2: Final partitions and cluster centroids obtained from the K-means algorithms with different initialisation. The best result is that of the middle figure.

the encoder function  $\text{enc}_\theta : \mathbb{R}^p \rightarrow \{1, \dots, K\}$  and decoder function  $\text{dec}_\theta : \{1, \dots, K\} \rightarrow \mathbb{R}^p$  as

$$\begin{aligned} z_i &= \text{enc}_\theta(x_i) = \underbrace{\arg \min_{k=1, \dots, K} \|x_i - \mu_k\|^2}_{\widehat{x}_i = \text{dec}_\theta(z_i) = \mu_{z_i}} \\ \widehat{x}_i &= \text{dec}_\theta(z_i) = \mu_{z_i} \end{aligned}$$

and let  $h_\theta(x_i) = \text{dec}_\theta(\text{enc}_\theta(x_i))$  be the corresponding autoencoder. For a squared loss function  $L$ , the risk is defined as

$$R(h_\theta) = \mathbb{E}[L(X, h_\theta(X))] = \mathbb{E}[\|X - h_\theta(X)\|^2]$$

The empirical risk is defined as

$$\begin{aligned} \widehat{R}_n(h_\theta) &= \frac{1}{n} \sum_{i=1}^n \|x_i - h_\theta(x_i)\|^2 \\ &= \underbrace{\min_{z_1, \dots, z_n} \frac{1}{n} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2}_{\widehat{\theta} = \arg \min_{\theta} \widehat{R}_n(h_\theta)} \\ &= \min_{C_1, \dots, C_K} \frac{1}{n} W(C_1, \dots, C_K, \mu_1, \dots, \mu_K) \end{aligned}$$

The value  $\widehat{\theta} = (\widehat{\mu}_1, \dots, \widehat{\mu}_K)$  that minimises  $W$  can be therefore be interpreted as the empirical risk minimiser

$$\widehat{\theta} = \underbrace{\arg \min_{\theta} \widehat{R}_n(h_\theta)}_{\widehat{\theta}}.$$

### 3.2.4 Example: Crabs dataset

We revisit the crabs dataset discussed in the chapter on PCA. Recall that each data  $x_i \in \mathbb{R}^5$ . We wish to apply K-means on the unlabelled crabs data, and use the class information (species and sex), as ground truth to check the partition obtained. We first apply K-means on the raw data. The results are given in Figure 3.3. The two clusters found by K-means divide the crabs between small crabs and large crabs, which is not really what we aim to find. As seen in the chapter on PCA, the size of the crabs explains most of the variance, and this is captured by the first PC. The projections onto the second and third components capture information about the different classes. We will therefore apply the K-means algorithm to the projections of the data onto those components, rescaled by the eigenvalues to have unit variance. For  $i = 1, \dots, n$ ,  $j = 2, 3$ , let

$$\tilde{z}_{ij} = \frac{v_j^\top x_i}{\sqrt{\lambda_j}}$$

where  $v_j$  is the  $j$ th PC and  $\lambda_j$  is the corresponding  $j$ th eigenvalue. We apply K-means on  $((\tilde{z}_{12}, \tilde{z}_{13})^\top, \dots, (\tilde{z}_{n2}, \tilde{z}_{n3})^\top)$ . The results are shown in Figure 3.4. For  $K = 2$ , the obtained partition correspond to a male/female partition of the crabs. For  $K = 4$ , we recover the four different classes BM, BF, OM and OF.

### 3.2.5 Additional comments

A good practice for initialisation is to randomly pick  $K$  training examples (without replacement) and set  $\mu_1, \mu_2, \dots, \mu_K$  equal to those examples. The Euclidean distance can be greatly affected by measurement unit and by strong correlations. Instead of the Euclidian distance, we can use the Mahalanobis distance instead:

$$\|x - y\|_M = \sqrt{(x - y)^\top M^{-1}(x - y)}$$

where  $M$  is a positive semi-definite matrix, e.g. the sample covariance. The partition obtained by the K-means algorithm will be very sensitive to outliers, as those outliers will drive the cluster centroids far from the other data examples in the cluster. The K-means algorithm works well when the clusters have a spherical shape, but will work poorly with non-convex cluster shapes.

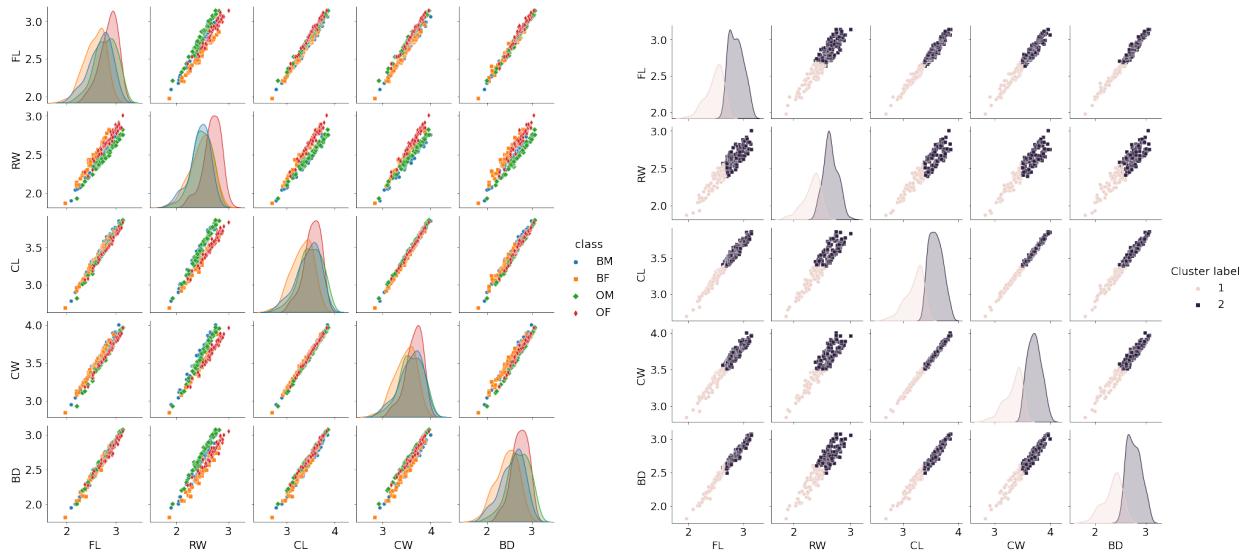


Figure 3.3: Left: pairplot of the crabs data, with the class information. Right: result of the K-means algorithm on the unlabelled data, with  $K = 2$  clusters.

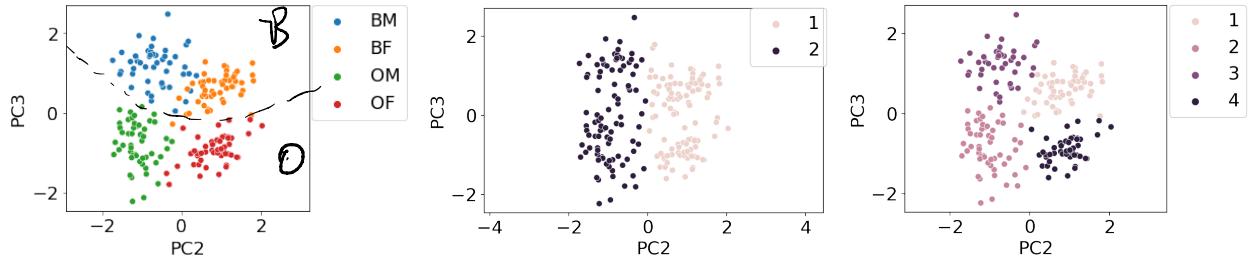


Figure 3.4: Left: Projection of the Crabs data onto the second and third PC, with the class information. Middle: Result of the K-means algorithm with  $K = 2$  clusters. Right: Result of the K-means algorithm with  $K = 4$  clusters.

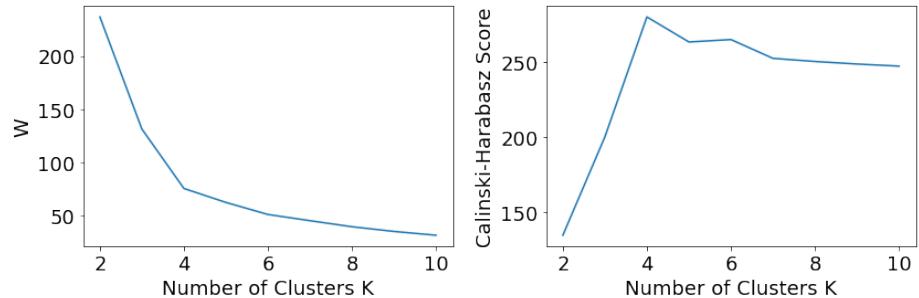


Figure 3.5: Left: Value of the objective function  $W$  after convergence for different number of clusters  $K$  on the Crabs dataset. Right: Calinski-Harabasz score as a function of the number  $K$  of clusters for the Crabs dataset.

### 3.2.6 Choice of $K$

The K-means objective function  $W$  will always improve with larger number of clusters  $K$ . For instance, Figure 3.5 shows the value of the objective function  $W$  after convergence for different number of clusters on the Crabs dataset. The function typically first decreases quickly then decreases more slowly. A heuristic approach, known as the elbow, looks for the value  $K$  where the curve bends from a high slope to a low slope.

Alternatively, one can consider using a different criterion to choose  $K$ . For example, the Calinski-Harabasz score, defined as

$$CH = \frac{\sum_{k=1}^K |C_k| \times \|\mu_k - \bar{x}\|^2}{\sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2} \times \frac{n - K}{K - 1}$$

is a measure of the separability of the clusters: higher values correspond to dense and well-separated clusters. For example, on the Crabs dataset, the value  $K = 4$  maximises this score (see Figure 3.5 (right)).



### 3.2.7 Stochastic optimisation

Each iteration of K-means requires a pass through whole dataset. In extremely large datasets, this can be computationally prohibitive. An alternative to using coordinate descent to optimise  $W$  is to resort to stochastic optimisation: we update the cluster centroids after assigning each data point to the closest cluster. We will review stochastic optimisation methods in more details later in the course, and just state the algorithm for K-means here. Let  $(\alpha_t)$  be a sequence of step-sizes satisfying  $\sum_{t=1}^{\infty} \alpha_t = \infty$  and  $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$ .

---

#### Algorithm 3 Stochastic Optimisation for K-means

---

- Repeat for  $t = 1, 2, \dots$  until satisfactory convergence:
  1. Pick data item  $x_i$  either randomly or in order.
  2. Assign  $x_i$  to the cluster with the nearest centroid,

$$z_i := \arg \min_{k=1, \dots, K} \|x_i - \mu_k\|^2$$

3. Update cluster centroid:

$$\mu_{z_i} := \mu_{z_i} + \alpha_t (x_i - \mu_{z_i}).$$


---

### 3.2.8 K-means for vector quantisation

A related algorithm, known as vector quantisation, has been developed in the signal processing literature for *lossy data compression*. The data matrix  $\mathbf{X}$  is represented by  $n \times p$  real numbers. Instead of the  $n \times p$  number, vector quantisation stores instead

- the *codebook* of  $K$  codewords  $\theta = (\mu_1, \dots, \mu_K)$  ( $K \times p$  real numbers)
- for each vector  $x_i$  its cluster assignment  $z_i$  ( $\lceil \log K \rceil \times n$  bits).

Each example  $x_i$  example is then approximated by the reconstructed example  $\hat{x}_i = \text{dec}_{\theta}(z_i)$ . As with K-means,  $K$  must be specified. Increasing  $K$  improves the quality of the compressed image but worsens the data compression rate, so there is a clear tradeoff. Some audio and video codecs use this method. The stochastic optimization algorithm for K-means was originally developed for vector quantisation.

**Example 9.** For example, consider  $3 \times 3$  block vector quantisation: we view each block of  $3 \times 3$  pixels of an image in RGB colors as a single observation  $x_i \in \mathbb{R}^{27}$ , called a super-pixel. An image is composed of  $n$  superpixels  $(x_1, \dots, x_n)$ . We apply K-means on the superpixels, and store the cluster labels and centroids instead of the raw data. The original images, and the decoded images for different values of  $K$  are represented in Figure 3.6.



Figure 3.6: Illustration of vector quantisation with various numbers of clusters  $K$  (codebook length).

---

SB2.2/SM4 Statistical Machine Learning  
Lecture notes - Supervised Learning

---

François Caron

University of Oxford, Hilary Term 2021

Version of January 24, 2021  
Please report typos to [caron@stats.ox.ac.uk](mailto:caron@stats.ox.ac.uk).

# Aims and Structure of this course

The aims of this course are

- To learn a number of different machine learning methods and gain understanding about their statistical foundations;
- To learn to identify and use the appropriate method for a given dataset and a given task;
- To learn how to use the relevant Python libraries/modules to analyse data, interpret results and evaluate the methods.

The course is structured in two main parts. In the first part, we will cover unsupervised learning (dimensionality reduction, feature extraction and clustering). In the second part of the course, we will cover supervised learning (classification and regression).

## References

These lecture notes are partially based on the slides from previous lecturers who have taught this course, and on the following books:

- L. Wasserman. All of Statistics. Springer, 2010.
- K. Murphy. Machine Learning. A probabilistic perspective. The MIT Press, 2012
- C. M. Bishop. Pattern Recognition and Machine Learning. Springer, 2006.
- T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical Learning. Springer, 2009.
- I. Goodfellow, Y. Bengio and A. Courville. Deep Learning. The MIT Press, 2016.

## Software

This document includes a number of examples written in the programming language Python. Python is certainly the most popular programming language for machine learning; it can be freely downloaded and installed. If you are new to Python, the following website contains useful information for beginners: <https://www.python.org/about/gettingstarted/>.

A number of Python code editors are freely available: PyCharm, Spyder, Atom, Jupyter. You may want to consider installing Anaconda, which is a distribution of Python and R programming languages for data science and machine learning. It comes with code editors such as Spyder and Jupyter for Python.

Python contains a very rich collection of libraries/packages for data science/machine learning. We will use the following libraries in this course:

- Numpy and Scipy: Core toolboxes for scientific computing (mathematical functions, linear algebra routines, random number generators, optimisation, etc.)
- Matplotlib: Visualisation
- Pandas: Data manipulation
- Seaborn: Statistical visualisation
- Scikitlearn: Machine Learning
- Pytorch: Deep Learning

These libraries can be loaded with the following commands.

```
import numpy as np
import scipy as sp
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import sklearn as skl
import torch
```

## Background material

### 0.1 Notations

Let  $x_1, \dots, x_p$  be scalar values. A  $p$ -dimensional row vector is noted as  $(x_1, \dots, x_p)$ . A  $p$ -dimensional column vector is noted as  $(x_1, \dots, x_p)^\top$ . By default, a vector means a column vector. The L2 norm of a vector  $x = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$  is defined as

$$\|x\| = \sqrt{x^\top x} = \sqrt{\sum_{j=1}^p x_j^2}$$

Let  $X \in \mathcal{X} \subseteq \mathbb{R}^p$  be a random variable with cumulative distribution function  $F$ . Throughout this course, we will use the notation

$$\mathbb{E}[h(X)] = \int_{\mathcal{X}} \phi(x) dF(x) = \begin{cases} \int_{\mathcal{X}} \phi(x) f(x) dx & \text{if } X \text{ is a continuous random variable} \\ \sum_{x \in \mathcal{X}} \phi(x) f(x) & \text{if } X \text{ is a discrete random variable} \end{cases}$$

where  $f$  denotes the probability density function of  $X$  if  $X$  is a continuous random variable, or the probability mass function of  $X$  if  $X$  is a discrete random variable.

**Definition 1** (Statistical Functional). *Let  $F$  be a cumulative distribution function. A **statistical functional** is a map  $T$  that maps a cdf  $F$  to a real number (or vector)  $T(F)$ .*

For example, for a random variable  $X$  with distribution  $F$ , its median  $m$  and mean  $\mu$  are both statistical functionals of  $F$ , with

$$m = F^{-1}(1/2), \quad \mu = \int_0^\infty x dF(x).$$

Let  $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$  be a random vector. We denote  $\mathbb{E}[X]$  the mean of  $X$  and

$$\text{cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top]$$

the  $p$ -by- $p$  covariance matrix of  $X$ . If  $p = 1$ , we write  $\text{var}(X)$  for its variance. For a vector  $a \in \mathbb{R}^p$ ,

$$\begin{aligned} \mathbb{E}[a^\top X] &= a^\top \mathbb{E}[X] \\ \text{var}(a^\top X) &= a^\top \text{cov}(X) a. \end{aligned}$$

The probability density function (pdf) of a Gaussian random vector  $X \in \mathbb{R}^p$  with mean  $\mu \in \mathbb{R}^p$  and  $p$ -by- $p$  covariance matrix  $\Sigma$  is

$$\varphi(x; \mu, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right). \quad (1)$$

### 0.2 Derivatives

Let  $J : \mathbb{R}^p \rightarrow \mathbb{R}$  be a scalar function. For a vector  $x = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ , the vector of partial derivatives, or gradient, is the  $p$ -dimensional vector given by

$$\nabla_x J(x) = \frac{\partial J(x)}{\partial x} = \begin{pmatrix} \frac{\partial J(x)}{\partial x_1} \\ \frac{\partial J(x)}{\partial x_2} \\ \vdots \\ \frac{\partial J(x)}{\partial x_p} \end{pmatrix}. \quad (2)$$

Let  $x, a \in \mathbb{R}^p$  and  $B$  be a  $p$ -by- $p$  symmetric matrix. We have

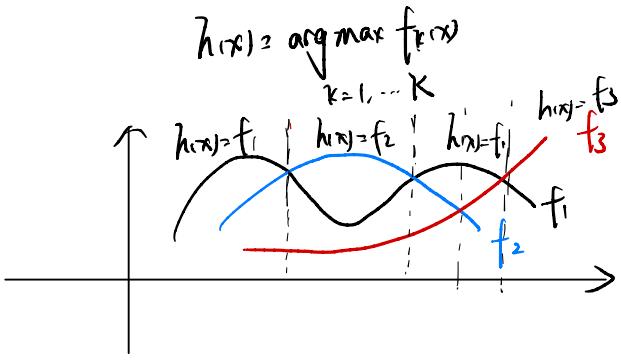
$$\frac{\partial(x^\top a)}{\partial x} = \frac{\partial(a^\top x)}{\partial x} = \underline{a} \quad (3)$$

$$\frac{\partial((x - a)^\top B(x - a))}{\partial x} = \frac{\partial((a - x)^\top B(a - x))}{\partial x} = \underline{2B(x - a)} \quad (4)$$

For a  $p$ -by- $p$  matrix  $\Sigma$ , let  $|\Sigma|$  denote the determinant of  $\Sigma$ . Let  $a, b \in \mathbb{R}^p$ . For a full-rank symmetric matrix  $\Sigma$ ,

$$\frac{\partial \log |\Sigma|}{\partial \Sigma} = \underline{\Sigma^{-1}} \quad (5)$$

$$\frac{\partial(a^\top \Sigma^{-1} b)}{\partial \Sigma} = \underline{-\Sigma^{-1} ab^\top \Sigma^{-1}} \quad (6)$$



# 1 — Basics of Statistical learning theory for supervised learning

## 1.1 Notations and assumptions

Let  $\mathcal{X}$  be the input space and  $\mathcal{Y}$  be the target space. We will focus here on the cases where  $\mathcal{Y} = \mathbb{R}$ , which corresponds to a regression task, or  $\mathcal{Y}$  is a finite set with  $K \geq 2$  elements, which corresponds to a classification task. By convention, we will assume that  $\mathcal{Y} = \{1, \dots, K\}$  if  $K \geq 3$  (multiclass classification) and  $\mathcal{Y} = \{1, -1\}$  if  $K = 2$  (binary classification).

**Definition 2** (Prediction rule). Let  $\mathcal{X}$  be the input space and  $\mathcal{Y}$  be the target space. A **prediction rule** is a function  $h : \mathcal{X} \rightarrow \mathcal{Y}$ . Let  $\mathcal{F} = \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$  the set of all prediction rules.

In the classification case, the prediction rule is called a **classification rule**, or a **classifier**. Without loss of generality, classifiers can be (non-uniquely) specified via a collection of real functions  $(f_k)_{k=1, \dots, K}$ :

$$\underbrace{h(x) = \arg \max_{k=1, \dots, K} f_k(x)} \quad (1.1)$$

where  $f_k : \mathcal{X} \rightarrow \mathbb{R}$ ,  $k = 1, \dots, K$  are called **discriminant functions**.  $f_k(x)$  may be interpreted as the level of confidence for choosing class  $k$  for an input  $x$ . The multivariate function  $\tilde{h} : x \rightarrow (f_1(x), \dots, f_K(x))$  is sometimes called a **soft classifier** (by contrast,  $h$  is then called a **hard classifier**). For two classes  $k$  and  $k'$ , the set of solutions to the equation  $f_k(x) = f_{k'}(x)$  is called the **decision boundary** between these two classes. For a binary classifier  $h$ , one can consider a single discriminant function  $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\underbrace{h(x) = \begin{cases} 1 & \text{if } f(x) \geq 0 \\ -1 & \text{if } f(x) < 0. \end{cases}} \quad (1.2)$$

The decision boundary between the two classes of the binary classifier is given by the solutions to the equation  $f(x) = 0$ .

**Definition 3** (Linear prediction rule). A **linear prediction rule** is such that

$$\underbrace{h(x) = \beta_0 + \beta^T x} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad (1.3)$$

for regression, where  $\beta_0 \in \mathbb{R}$ ,  $\beta = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$ , and

$$\underbrace{f_k(x) = \beta_{k0} + \beta_k^T x} \quad (1.4)$$

for classification, where  $\beta_{k0} \in \mathbb{R}$  and  $\beta_k = (\beta_{k1}, \dots, \beta_{kp}) \in \mathbb{R}^p$  for  $k = 1, \dots, K$ . The decision boundary between two classes  $k$  and  $k'$  is the hyperplane defined by the equation

$$\underbrace{(\beta_{k0} - \beta_{k'0}) + (\beta_k - \beta_{k'})^T x = 0.} \quad \underbrace{f_k(x) - f_{k'}(x) = 0}$$

Consider a dataset  $d = (x_i, y_i)_{i=1, \dots, n}$  where  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$  is the  $i$ th input/target observation.  $n \geq 1$  is the number of observations. Our objective is, based on the dataset  $d$ , to learn (or train) a prediction rule  $\hat{h}^{(d)} \in \mathcal{F}$ . This prediction rule is then used, for a set of new examples with known input values  $(x_{n+1}, x_{n+2}, \dots)$ , to predict their target values  $(\hat{y}_{n+1}, \hat{y}_{n+2}, \dots)$  using the **learned prediction rule**

$$\underbrace{\hat{y}_{n+i} = \hat{h}^{(d)}(x_{n+i})}.$$

We give below two simple examples of learned prediction rules  $\hat{h}^{(d)}$  for regression and classification.

**Example 4** (Least square linear regression). Consider univariate regression ( $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ ), with the least square prediction rule, defined as

$$\hat{h}^{(d)}(x) = \hat{\beta}_0 + x\hat{\beta}_1 \quad (1.5)$$

where the least square estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are defined as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (1.6)$$

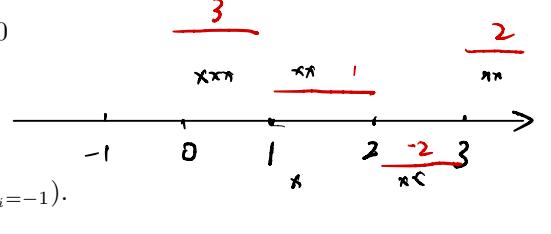
where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . The prediction rule  $\hat{h}^{(d)}(x)$  is linear.

**Example 5** (Histogram binary classifier). For  $\mathcal{X} = \mathbb{R}$  and  $\mathcal{Y} = \{-1, 1\}$  consider the histogram classifier

$$\hat{h}^{(d)}(x) = \begin{cases} 1 & \text{if } \hat{f}^{(d)}(x) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

where the discriminant function  $\hat{f}^{(d)}$  is defined as

$$\hat{f}^{(d)}(x) = \sum_{i=1}^n \mathbb{1}_{[x_i] = [x]} (\mathbb{1}_{y_i=1} - \mathbb{1}_{y_i=-1}).$$



This simple classifier is piecewise constant, and assigns a point  $x$  in the interval  $[x, x+1)$  to the most frequent class in that interval in the dataset  $d$ .

We assume that the dataset  $d = (x_i, y_i)_{i=1,\dots,n}$  is a realisation from a random sample  $D = (X_i, Y_i)_{i=1,\dots,n}$  where  $(X_1, Y_1), \dots, (X_n, Y_n)$  are independent and identically distributed (iid) random variables with the same distribution as  $(X, Y)$ . Let  $P_0$  denote the distribution of  $(X, Y)$ , called the true data distribution, or population distribution. Note that  $\hat{h}^{(D)}$  therefore denotes a random prediction rule, as it is parameterised by the random sample  $D$ ; the prediction rule  $\hat{h}^{(d)}$  is a realisation of  $\hat{h}^{(D)}$ .

**Decomposition of the variance.** For classification with  $K$  classes, using the law of total variance, we have the following decomposition of the covariance matrix of the input vector  $X$

$$\text{cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top] = \underbrace{\mathbb{E}[\text{cov}(X | Y)]}_{\text{within-class covariance}} + \underbrace{\text{cov}(\mathbb{E}[X | Y])}_{\text{between-class covariance}}. \quad (1.7)$$

The first term in the right-handside of (1.7) corresponds to the part of the total covariance matrix  $\text{cov}(X)$  originating from the variance **within** each class. The second term corresponds to the variance originating from the variance **between** the different classes. Let  $\mu_k = \mathbb{E}[X | Y = k]$ ,  $\Sigma_k = \text{cov}(X | Y = k)$ ,  $\mu = \mathbb{E}[X]$ , and  $\pi_k = \Pr(Y = k)$ . The two covariance terms can be expressed as

$$B = \text{cov}(\mathbb{E}[X | Y]) = \sum_{k=1}^K \pi_k (\mu_k - \mu)(\mu_k - \mu)^\top$$

$$W = \mathbb{E}[\text{cov}(X | Y)] = \sum_{k=1}^K \pi_k \Sigma_k$$

$\check{E}(X) = \check{E}(\check{E}(X | Y=k))$

If  $X \in \mathbb{R}$ , the ratio

$$\frac{\text{cov}(\mathbb{E}[X | Y])}{\mathbb{E}[\text{cov}(X | Y)]} \geq 0$$

between the between-class and within-class variances is a measure of the **separability** between the classes. Larger values indicate more separated classes, while small values indicate a lot of overlap between the classes. For instance, if  $\Pr(Y = 1) = \Pr(Y = -1) = 1/2$  and  $X | Y = y \sim \mathcal{N}(y\mu, \sigma^2)$ , for some  $\mu \geq 0$  and  $\sigma > 0$ , we have

$$\mathbb{E}[\text{cov}(X | Y)] = \sigma^2, \quad \text{cov}(\mathbb{E}[X | Y]) = \mu^2.$$

$\sigma^2$  therefore controls the within-class variance, and  $\mu^2$  the between-class variance.

In order to formalise mathematically the supervised learning task, we need to answer the following questions:

- How do we quantify the accuracy of a prediction rule?
- If we were to know the true data distribution, what would be the associated (optimal) prediction rule?
- Based on a dataset of size  $n$ , how can we learn a prediction rule whose accuracy is close to the optimal one?

## Law of Total Variance:

### ① Law of Total Expectation:

$$\bar{E}(Y) = \bar{E}(\bar{E}(Y|X)) \Rightarrow \sum_i \bar{E}(Y|X_i) P(X_i)$$

$$② \text{Var}(Y) = \bar{E}(Y^2) - \bar{E}(Y)^2$$

Since  $\bar{E}(Y^2) > \bar{E}[\text{Var}(Y|X) + \bar{E}(Y|X)^2]$

$$\begin{aligned} \Rightarrow \text{Var}(Y) &= \bar{E}[\text{Var}(Y|X) + \bar{E}(Y|X)^2] - \bar{E}[\bar{E}(Y|X)]^2 \\ &= \bar{E}[\text{Var}(Y|X)] + \underbrace{\bar{E}[\bar{E}(Y|X)^2] - \bar{E}[\bar{E}(Y|X)]^2}_{\text{Var}[\bar{E}(Y|X)]} \\ &= \bar{E}[\text{Var}(Y|X)] + \text{Var}[\bar{E}(Y|X)]. \end{aligned}$$

Similarly  $\text{Cov}(X) = \bar{E}[\text{Cov}(X|Y)] + \text{Cov}[\bar{E}(X|Y)]$

In order to formalise mathematically the supervised learning task, we need to answer the following questions:

How do we quantify the accuracy of a prediction rule?

Loss function and risk

If we were to know the true data distribution, what would be the associated (optimal) prediction rule?

Bayes prediction rule

Based on a dataset of size  $n$ , how to choose a prediction rule whose accuracy is close to the optimal one?

Empirical risk minimisation

$$\bar{E}(Y^2|X)$$

## 1.2 Loss function and risk

Consider a prediction rule  $h : \mathcal{X} \rightarrow \mathcal{Y}$  where  $h(x) \in \mathcal{Y}$  is the predicted target for an input  $x$ . We need a way to measure how “bad” it is to predict  $h(x)$  when the true target is indeed  $y$ . This will be done via the introduction of a loss function.

**Definition 6.** Let  $\mathcal{Y}$  be the target space. A **loss function** is a non-negative function

$$L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

The loss  $L(y, h(x))$  represents the “cost” of predicting a target  $h(x)$  when the true target is  $y$ . We typically have  $L(y, h(x)) = 0$  when  $h(x) = y$ , and  $L(y, h(x))$  increases as  $y$  and  $h(x)$  are “further away”. Typical loss functions are the squared error loss

$$\underline{L(y, h(x)) = (y - h(x))^2} \quad (1.8)$$

for regression, and the 0 – 1 loss for classification

$$\underline{L(y, h(x)) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}} \quad (1.9)$$

The 0 – 1 loss assigns a cost of 1 if the predicted class is different from the true target, and 0 otherwise.

**Definition 7 (Risk - supervised learning).** For a given loss function  $L$ , the (population) **risk**  $R(h)$  of a prediction rule  $h$  is the expected loss

$$R(h) = \mathbb{E}[L(Y, h(X))]$$

where the expectation is taken over the true (unknown) joint distribution of  $(X, Y)$ .

The risk is therefore given by

$$R(h) = \begin{cases} \mathbb{E}[(Y - h(X))^2] & \text{under the squared loss (1.8)} \\ \Pr(Y \neq h(X)) & \text{under the 0 - 1 loss (1.9).} \end{cases}$$

$R(h)$  is therefore the mean squared error under the squared loss, and the probability of misclassification for classification under the 0-1 loss. We assume for simplicity that  $R(h)$  admits a unique minimum over the set of prediction rules  $\mathcal{F}$ . The optimal prediction rule, called **Bayes prediction rule**, is the function  $h^*$  that minimises the risk  $R(h)$ .

**Definition 8** (Bayes prediction rule). The Bayes prediction rule is given by

$$h^* = \arg \min_{h \in \mathcal{F}} R(h). \quad (1.10)$$

 **Proposition 9.** For any  $x \in \mathcal{X}$ , we have

$$h^*(x) = \arg \min_{y^* \in \mathcal{Y}} \mathbb{E}[L(Y, y^*) | X = x], \quad (1.11)$$

where the expectation is taken over the true (unknown) conditional distribution of  $Y$  given  $X = x$ .

*Proof.* Note that for any function  $f$ ,  $\mathbb{E}[f(X, Y)] = \mathbb{E}_X[\mathbb{E}_Y[f(X, Y) | X]]$ . It follows that, for any prediction rule  $h$ ,

$$\underline{R(h) = \mathbb{E}_X[\mathbb{E}_Y[L(Y, h(X)) | X]]}, \quad R(h) = \mathbb{E}[L(Y, h(X))]$$

As, for any  $x$ ,

$$\begin{aligned} \mathbb{E}_Y[L(Y, h(X)) | X = x] &\geq \min_{y \in \mathcal{Y}} \mathbb{E}_Y[L(Y, y) | X = x] \\ &= \mathbb{E}_Y[L(Y, h^*(X)) | X = x] \end{aligned}$$

we therefore obtain, for any  $h \in \mathcal{F}$ ,  $\underline{R(h) \geq R(h^*)}$ . □

Let  $F_x(y) = \Pr(Y \leq y \mid X = x)$  denote the conditional cdf of  $Y$  given  $X = x$ . For each  $x$ , the Bayes prediction  $h^*(x)$  is a statistical functional of the conditional cdf  $F_x$ :

$$\underline{h^*(x) = T(F_x)} \quad (1.12)$$

with

$$\underline{T(F_x) = \arg \min_{y^* \in \mathcal{Y}} \int_{\mathcal{Y}} L(y, y^*) dF_x(y)}. \quad (1.13)$$

The risk  $R(h^*)$  of the Bayes prediction function  $h^*$  is called the **Bayes risk**;  $h^*$  is optimal in the sense that it achieves the lowest risk amongst all prediction rules. That is, for any function  $h \in \mathcal{F}$ , we have

$$\underline{R(h) \geq R(h^*)}.$$

The difference

$$\underline{R(h) - R(h^*) \geq 0} \quad (1.14)$$

is called the **excess risk** of the prediction rule  $h$ .

**D** In the case of the squared loss function and of the  $0 - 1$  loss, the Bayes prediction rule and the Bayes risk both admit a simple expression. For a squared loss function, we obtain (see Problem Sheet)

$$\left\{ \begin{array}{l} h^*(x) = \mathbb{E}[Y \mid X = x] \\ R(h^*) = \mathbb{E}_X[\text{var}(Y \mid X)] \\ R(h) - R(h^*) = \mathbb{E}_X[(h(X) - h^*(X))^2]. \end{array} \right.$$

In the case of the  $0 - 1$  loss, we obtain for classification problems

$$\left\{ \begin{array}{l} h^*(x) = \arg \max_{k \in \mathcal{Y}} \Pr(Y = k \mid X = x) \\ R(h^*) = 1 - \mathbb{E}_X \left[ \max_{k \in \mathcal{Y}} \Pr(Y = k \mid X) \right]. \end{array} \right.$$

#### Bayes prediction rule

For example, let  $(X, Y)$  be defined as

$$X \sim \mathbf{U}(0, 1)$$

and given  $X = x$ ,

$$Y = m(x) + \varepsilon$$

where  $\varepsilon$  is independent of  $X$ , with  $\mathbb{E}[\varepsilon] = 0$ ,  $\text{var}(\varepsilon) = \sigma^2$ .

The Bayes prediction rule is

$$h^*(x) = \mathbb{E}[Y \mid X = x] = \mathbb{E}[m(X) + \varepsilon \mid X = x] = m(x)$$

with Bayes risk

$$R(h^*) = \mathbb{E}[\text{var}(Y \mid X)] = \sigma^2$$

The risk of a prediction rule  $h$  is

$$\begin{aligned} R(h) &= \mathbb{E}[(Y - h(X))^2] = \mathbb{E}[(m(X) + \varepsilon - h(X))^2] \\ &= \mathbb{E}[(m(X) - h(X))^2] + \mathbb{E}[\varepsilon^2] = \underbrace{\int_0^1 (m(x) - h(x))^2 dx}_{\text{Excess risk}} + \sigma^2 \end{aligned}$$

The discriminant functions of the optimal classifier are therefore  $f_k(x) = \Pr(Y = k \mid X = x)$ . For binary classification, this can be further simplified to

$$\left\{ \begin{array}{l} h^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq 1/2 \\ -1 & \text{otherwise} \end{cases} \\ R(h^*) = \frac{1}{2} - \mathbb{E}[|\eta(X) - 1/2|] \end{array} \right. \quad (1.20)$$

where  $\eta(x) = \Pr(Y = 1 \mid X = x)$ . The associated discrimination function is therefore  $f(x) = \eta(x) - \frac{1}{2}$ . Additionally, the excess risk takes the form

$$\underline{R(h) - R(h^*) = \mathbb{E}_X[(2\eta(X) - 1)(\mathbb{1}_{h^*(X)=1} - \mathbb{1}_{h(X)=1})]} \quad (1.22)$$

### 1.3 Optimal classifier under Gaussian class distributions

We consider here the optimal classifier under Gaussian class distributions. Assume that the joint distribution of  $(X, Y)$ , where  $X \in \mathbb{R}^p$  and  $Y \in \{1, \dots, K\}$ , is defined as follows.

$$\Pr(Y = k) = \pi_k \quad (1.23)$$

$$\underline{X \mid Y = k \sim \mathcal{N}(\mu_k, \Sigma_k)} \quad (1.24)$$

where  $\pi_k \geq 0$  are the class probabilities, with  $\sum_{k=1}^K \pi_k = 1$ ,  $\mu_k \in \mathbb{R}^p$  is the mean of observations in class  $k$  and  $\Sigma_k$  is the  $p \times p$  covariance matrix of observations in class  $k$ . Denote

$$g_k(x) = (2\pi)^{-p/2} |\Sigma_k|^{-1/2} \exp \left( -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right) \quad (1.25)$$

the conditional pdf of  $X$  given  $Y = k$ . An application of Bayes' theorem gives

$$\Pr(Y = k \mid X = x) = \frac{\pi_k g_k(x)}{\sum_{\ell=1}^K \pi_\ell g_\ell(x)}.$$

$$\mathbb{P}(Y=k \mid X=x) = \sum \mathbb{P}_\ell(X=x \mid Y=\ell) \mathbb{P}^\ell(Y)$$

$$\mathbb{P}(Y=k \mid X=x) = \mathbb{P}(X=x \mid Y=k) \overline{\mathbb{P}(Y=k)} / \mathbb{P}(X=x)$$

Note that the denominator in Equation (2.1) does not depend on  $k$ . It follows that, under a  $0 - 1$  loss, the optimal classifier is

$$\begin{aligned} h^*(x) &= \arg \max_{k \in \{1, \dots, K\}} \Pr(Y = k | X = x) \\ &= \arg \max_{k \in \{1, \dots, K\}} \pi_k g_k(x) \\ &= \arg \max_{k \in \{1, \dots, K\}} \log(\pi_k) + \log(g_k(x)) \\ &= \arg \max_{k \in \{1, \dots, K\}} f_k(x) \end{aligned}$$

where the discriminant functions  $f_k$ ,  $k = 1, \dots, K$ , are defined as

$$\begin{aligned} f_k(x) &= \log(\pi_k) - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^\top \Sigma_k^{-1} (x - \mu_k) \\ &= \underbrace{\log(\pi_k) - \frac{1}{2} (\log |\Sigma_k| + \mu_k^\top \Sigma_k^{-1} \mu_k)}_{a_k} + \underbrace{\mu_k^\top \Sigma_k^{-1} x}_{b_k^\top x} - \underbrace{\frac{1}{2} x^\top \Sigma_k^{-1} x}_{c_k x^\top c_k x} \end{aligned}$$

where  $a_k = \log(\pi_k) - \frac{1}{2} (\log |\Sigma_k| + \mu_k^\top \Sigma_k^{-1} \mu_k)$ ,  $b_k = \Sigma_k^{-1} \mu_k$  and  $c_k = -\frac{1}{2} \Sigma_k^{-1}$ . The discriminant function is therefore a quadratic function of  $x$ . The decision boundary between class  $k$  and class  $k'$  is obtained by looking at the solutions to the equation  $f_k(x) - f_{k'}(x) = 0$ . This gives

$$a_k + b_k^\top x + x^\top c_k x - (a_{k'} + b_{k'}^\top x + x^\top c_{k'} x) = \underbrace{a_* + b_*^\top x + x^\top c_* x}_{} = 0.$$

where  $a_* = a_k - a_{k'}$ ,  $b_* = b_k - b_{k'}$  and  $c_* = c_k - c_{k'}$ . The decision boundary is therefore given by the roots of a quadratic function of  $x$ .

### Equal covariance and linear classifier.

If the covariance matrices are all equal,  $\Sigma_k = \Sigma$  for all  $k$ , the discriminant functions simplify to

$$\begin{aligned} f_k(x) &= \log(\pi_k) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} \underbrace{(x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)}_{\text{Mahalanobis distance}} \\ &= -\frac{1}{2} (x^\top \Sigma^{-1} x + \log |\Sigma|) + \underbrace{\log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k}_{\text{constant}} + \underbrace{\mu_k^\top \Sigma^{-1} x}_{\text{linear term}} \end{aligned}$$

As the first term  $-\frac{1}{2} (x^\top \Sigma^{-1} x + \log |\Sigma|)$  does not depend on  $k$ ,

$$\arg \max_{k \in \{1, \dots, K\}} f_k(x) = \arg \max_{k \in \{1, \dots, K\}} \tilde{f}_k(x)$$

with discriminant functions

$$\begin{aligned} \tilde{f}_k(x) &= \log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \mu_k^\top \Sigma^{-1} x \\ &= a_k + b_k^\top x \end{aligned}$$

where  $a_k = \log(\pi_k) - \frac{1}{2} (\mu_k^\top \Sigma^{-1} \mu_k)$ ,  $b_k = \Sigma^{-1} \mu_k$ . The discriminant functions are therefore linear, and the decision boundary between two classes is an hyperplane, solution to the equation

$$a_* + b_*^\top x = 0$$

where  $a_* = a_k - a_{k'}$  and  $b_* = b_k - b_{k'}$ . The resulting optimal classifier is therefore a **linear classifier**. Figure 1.1 provides illustrations of the Bayes classifier and decision boundaries when  $X$  is two-dimensional, and  $Y$  has two classes, when the covariance matrices are different or equal.

### Equal covariance and discriminant coordinates.

In the case of equal covariances, the predicted class  $h^*(x)$  only depends on  $x$  through lower-dimensional projections of  $x$ , called its **discriminant coordinates**. In order to gain some intuition, consider first the binary

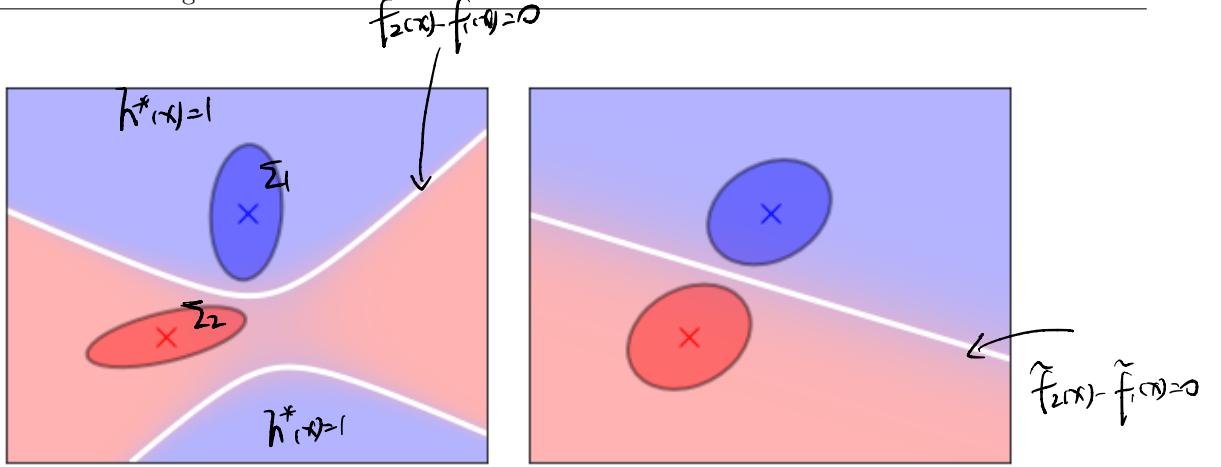


Figure 1.1: Optimal Bayes classifiers (background color) and associated decision boundaries (white line) under Gaussian class distributions, with (left) different covariances within each class and (right) the same covariance. Crosses correspond to the class means, and ellipses represent the contour of constant pdf of the Gaussian class distributions.

$$\tilde{f}_k(x) = \log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \frac{1}{2} \Sigma^{-1} x$$

case, with a discriminative function  $\eta(x) = \tilde{f}_2(x) - \tilde{f}_1(x)$

$$\begin{aligned} \eta(x) &= \log(\pi_2) - \log(\pi_1) - \frac{1}{2} (\mu_2^\top \Sigma^{-1} \mu_2 - \mu_1^\top \Sigma^{-1} \mu_1) + (\mu_2 - \mu_1)^\top \Sigma^{-1} x \\ &= \log \frac{\pi_2}{\pi_1} - \frac{1}{2} (\mu_2 - \mu_1)^\top \Sigma^{-1} (\mu_2 + \mu_1) + (\mu_2 - \mu_1)^\top \Sigma^{-1} x \end{aligned}$$

Letting

$$x_0 = \frac{1}{2} (\mu_1 + \mu_2) - (\mu_2 - \mu_1) \frac{\log(\pi_2/\pi_1)}{(\mu_2 - \mu_1)^\top \Sigma^{-1} (\mu_2 - \mu_1)}$$

we obtain

$$\eta(x) = (\mu_2 - \mu_1)^\top \Sigma^{-1} (x - x_0).$$

Note that  $x_0$  lies on the line passing through the class means  $\mu_1$  and  $\mu_2$ , and  $x_0 = \frac{1}{2}(\mu_1 + \mu_2)$  is half-way if  $\pi_2 = \pi_1 = 1/2$ .

For a vector  $x \in \mathbb{R}^p$ , writing  $x^\bullet = \Sigma^{-1/2}x$ , we obtain

$$\eta(x) = (\mu_2^\bullet - \mu_1^\bullet)^\top (x^\bullet - x_0^\bullet).$$

$$\mu_2^\bullet = \Sigma^{-1/2} \mu_2 \quad \mu_1^\bullet = \Sigma^{-1/2} \mu_1$$

The decision boundary therefore only depends on  $x$  through the projection of  $x^\bullet = \Sigma^{-1/2}x$  onto the vector of the difference of the transformed class means  $(\mu_2^\bullet - \mu_1^\bullet)$ . Hence, for the purpose of classification, only the one-dimensional projection  $(\mu_2^\bullet - \mu_1^\bullet)\Sigma^{-1/2}x$  of the  $p$ -dimensional vector  $x$  is relevant.

For general  $K$ , the classifier relies on the Mahalanobis distances

$$(x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)$$

between the vector  $x$  to the class means or, equivalently, on the squared Euclidian distances

$$(x^\bullet - \mu_k^\bullet)^\top (x^\bullet - \mu_k^\bullet)$$

between the transformed  $x^\bullet$  and  $\mu_k^\bullet$ . The  $K$  spheroid class means  $\mu_k^\bullet$  lie in an affine manifold of dimension lower or equal to  $K-1$ . When looking at the closest class mean, one can ignore distances orthogonal to this subspace, as they contribute equally for each class. We can therefore project the  $x^\bullet$  onto this class-spanning subspace  $H_{K-1}$ . This corresponds to finding the principal components subspace of the transformed means, or similarly, to find the eigendecomposition of the spheroid between-class covariance

$$B^\bullet = \text{cov}(\mathbb{E}[\Sigma^{-1/2}X|Y]) = \Sigma^{-1} \text{cov}(\mathbb{E}[X|Y]) = \sum_{k=1}^K \pi_k (\mu_k^\bullet - \mu^\bullet) (\mu_k^\bullet - \mu^\bullet)^\top$$

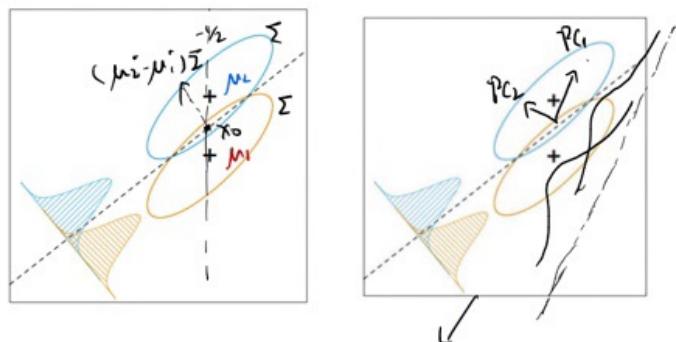
where  $\mu^\bullet = \sum_{k=1}^K \pi_k \mu_k^\bullet$ . Denote

$$B^\bullet = U^\bullet D^\bullet U^{\bullet\top}$$

## Equal covariances: discriminant coordinates

If we ignore PCA

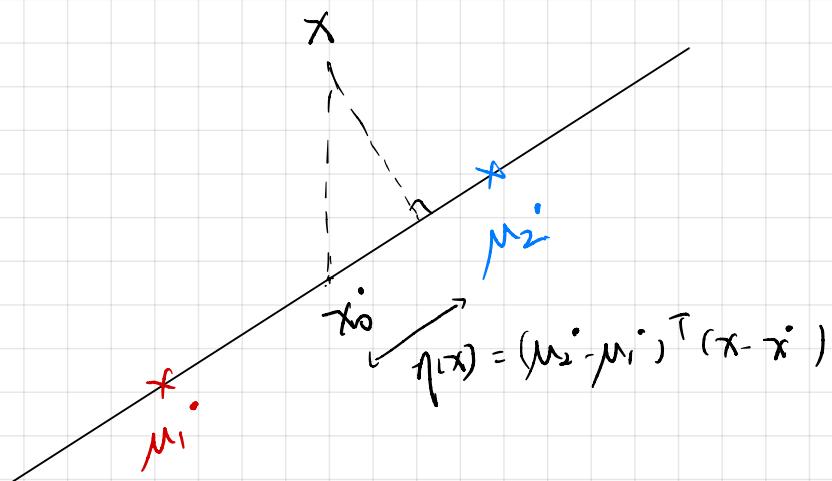
$$\pi_1 = \pi_2$$



Maximize Variance  
not best separate the class

18

Figure from Hastie, Tibshirani and Friedman, Section 4.3.3



$$\text{For general } k, \quad f_k(x) = \log(\pi_k) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)$$

Mahalanobis distance

The classifier relies on the Mahalanobis distance

$$(x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)$$

between the vector  $x$  to the class mean, or equivalently on the squared Euclidean distance

$$(x - \mu_k)^\top (x - \mu_k) = \|x - \mu_k\|^2$$

between the transformed  $x'$  and  $\mu_k'$ .

the eigendecomposition of  $B^\bullet$ . The columns  $u_\ell^\bullet$  of  $U^\bullet$  in sequence define the coordinates of the optimal subspaces. We therefore define

$$z_\ell = (u_\ell^\bullet)^\top x^\bullet = (u_\ell^\bullet)^\top \Sigma^{-1/2} x$$

the  $\ell$ th **discriminant coordinate** of the vector  $x$ .

Fisher arrived at the same decomposition with a different route, without referring to Gaussian distributions. Assume we want to find a linear combination  $Z = a^\top X$  such that the between-class variance is maximised relative to the within-class variance. That is, we want to maximise

$$\frac{\text{var}(\mathbb{E}[a^\top X|Y])}{\mathbb{E}[\text{var}(a^\top X|Y)]} = \frac{a^\top B a}{a^\top \Sigma a}.$$

Setting  $u = \Sigma^{1/2}a$  yields

$$\frac{u^\top B^\bullet u}{u^\top u}.$$

Maximisation over  $u$  is achieved by the first eigenvector of  $B^\bullet$ , which corresponds to  $u_1^\bullet$  as above. The next eigenvector  $u_2^\bullet$  is obtained by finding the vector orthogonal to  $u_1^\bullet$  that maximises  $\frac{u_2^\top B^\bullet u_2}{u_2^\top u_2}$ , etc.

## 1.4 Statistical Learning Approaches

For a prediction rule  $h$  and a dataset  $d = (x_i, y_i)_{i=1,\dots,n}$ , the **empirical risk**, or **training error**, is defined as

$$\widehat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n L(y_i, h(x_i)). \quad (1.27)$$

This represents the average error over the dataset  $d$  used to train/learn the classifier. The population risk  $R(h)$ , defined in Definition 7, is also called the **generalisation error** or **out-of-sample** error. It represents the average error for a new observation out of the training set  $d$ , and is the quantity we aim to minimise. The **generalisation gap** is defined as the difference between the generalisation error and the training error

$$R(h) - \widehat{R}_n(h).$$

Our objective is to find a prediction rule  $h$  with small generalisation error  $R(h)$ . The conditional distribution of  $Y|X = x$  being unknown, we cannot calculate the optimal prediction rule  $h^*$ . For a given function  $h$ , we cannot either compute its generalisation error  $R(h)$  as this involves an expectation with respect to the true unknown distribution of  $(X, Y)$ . All we have is a realisation  $d = (x_i, y_i)_{i=1,\dots,n}$  from a random sample  $D = (X_i, Y_i)_{i=1,\dots,n}$ .

Considering either the representation (1.10) or (1.11) for the optimal prediction rule, there are two broad families of methods that we might consider.

1. **Empirical Risk Minimisation** (ERM): Replace the population risk  $R(h)$  by the empirical risk  $\widehat{R}_n(h)$ , and minimise this function over some set of functions  $\mathcal{H} \subset \mathcal{F}$

$$\widehat{h}^{(d)} = \arg \min_{h \in \mathcal{H}} \widehat{R}_n(h). \quad (1.28)$$

2. **Plug-in methods:** Obtain some estimate  $\widehat{F}_x$  of the conditional distribution of  $Y$  given  $X = x$ , and set

$$\widehat{h}^{(d)}(x) = T(\widehat{F}_x) = \arg \min_{y^* \in \mathcal{Y}} \int_{\mathcal{Y}} L(y, y^*) d\widehat{F}_x(y). \quad (1.29)$$

The two different approaches address the prediction problem by considering tasks of increasing difficulty: ERM learns directly the map from  $\mathcal{X}$  to  $\mathcal{Y}$ , without trying to estimate the conditional distribution of  $Y|X = x$ . The conditional approach aims at estimating the conditional distribution of  $Y|X = x$  in order to derive the prediction rule. Plug-in methods can themselves be separated between the conditional approach and the generative approach:

- 2a. **The conditional plug-in** method estimates directly the conditional distribution of  $Y$  given  $X = x$
- 2b. **The generative plug-in** method first estimates the joint distribution of  $(X, Y)$ , then derives an estimate of the conditional distribution of  $Y$  given  $X = x$  via Bayes' theorem.

While the conditional approach leaves the marginal distribution of  $X$  unspecified, the generative approach aims at estimating the joint distribution, and hence tackles a more difficult problem.

The ERM and conditional plug-in approaches are commonly referred as **discriminative learning**, as one tries to directly learn the relation from  $X$  to  $Y$ . The generative plug-in approach is known as **generative learning**, as one tries to learn the full joint distribution of  $(X, Y)$ , hence tries to learn the full generative model of the data.

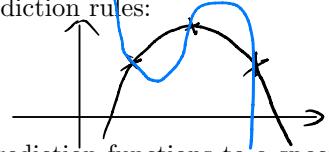
Method	ERM	Conditional plug-in	Generative plug-in
Learns...	Map from $X$ to $Y$	Conditional $Y X = x$	Joint dist. of $(X, Y)$
Type	Discriminative	Discriminative	Generative
Examples	Neural Nets Boosting Nearest neighbours Random Forests	Logistic regression	LDA/QDA Naive Bayes

### 1.4.1 Empirical Risk Minimisation: Definition

As the true data distribution of the data is unknown (as well as the conditional distribution of  $Y|X = x$ ), neither the true risk nor the optimal prediction function can be computed. However, we have a training set of data  $(x_i, y_i)_{i=1, \dots, n}$  which are realisations from the true data distribution  $P_0$ . As a proxy for the true risk, we use the empirical risk, which can be seen as a Monte Carlo estimate of the true risk  $R(h)$  as  $(x_i, y_i)_{i=1, \dots, n}$  are realisations of an iid random sample  $(X_i, Y_i)_{i=1, \dots, n}$  from  $P_0$ . As such, it is an unbiased and consistent estimator of the true risk.<sup>1</sup> For any fixed  $h$ ,  $\mathbb{E}[\widehat{R}_n(h)] = R(h)$  and  $\widehat{R}_n(h) \rightarrow R(h)$  almost surely as  $n \rightarrow \infty$  (See Problem Sheet).

We cannot minimise the empirical risk over the whole set of functions  $\mathcal{F}$ . If all the  $x_i$  take different values, the minimal empirical risk  $\widehat{R}_n(h) = 0$  is achieved for an infinite number of prediction rules:

$$\widehat{h}^{(d)}(x) = \begin{cases} y_i & \text{if } x \in \{x_1, \dots, x_n\} \\ h_0(x) & \text{otherwise} \end{cases}$$



where  $h_0 : \mathcal{X} \rightarrow \mathcal{Y}$  is an arbitrary function. We therefore restrict the set of prediction functions to a specified function space  $\mathcal{H} \subset \mathcal{F}$  (e.g. the set of linear classifiers).

**Definition 10.** For a function space  $\mathcal{H}$ , a loss function  $L$ , and a training set  $d = (x_i, y_i)_{i=1, \dots, n}$ , the Empirical Risk Minimiser  $\widehat{h}^{(d)}$  is defined as

$$\widehat{h}^{(d)} = \arg \min_{h \in \mathcal{H}} \widehat{R}_n(h) \quad (1.30)$$

where  $\widehat{R}_n(h)$  is the empirical risk defined in Equation (1.27).

In practice, a function  $h \in \mathcal{H}$  will often be parameterise by a vector  $\theta \in \Theta$ , and we write  $h_\theta$ . For instance, if  $\mathcal{H}$  is the class of linear functions on  $\mathbb{R}$ , then  $h_\theta(x) = \beta_0 + \beta_1 x$  where  $\theta = (\beta_0, \beta_1) \in \mathbb{R}^2$ . The ERM is therefore given by  $\widehat{h}^{(d)} = h_{\widehat{\theta}}$  where

$$\widehat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n L(y_i, h_\theta(x_i)).$$

### 1.4.2 Plug-in methods

Plug-in methods obtain some estimate  $\widehat{F}_x$  of the conditional distribution of  $Y$  given  $X = x$ , and set

$$\widehat{h}^{(d)}(x) = T(\widehat{F}_x) = \arg \min_{y^* \in \mathcal{Y}} \int_{\mathcal{Y}} L(y, y^*) d\widehat{F}_x(y). \quad (1.31)$$

We will focus here on the case where the estimate is obtaining via maximum likelihood. Assume that the conditional probability mass/density function of  $Y$  given  $X = x$  is  $f_\theta(y|x)$  where the parametric form  $f_\theta$  is known, but  $\theta \in \Theta$  is some unknown parameter. In the generative case,  $f_\theta(y|x) \propto \pi_\theta(x, y)$ , where  $\pi_\theta$  is the known parametric form of the joint distribution. The maximum likelihood estimate of  $\theta$  based on the training dataset  $d$  is

$$\left\{ \begin{array}{l} \widehat{\theta} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(f_\theta(y_i|x_i)) \quad [\text{Conditional Approach}] \\ \widehat{\theta} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(\pi_\theta(x_i, y_i)) \quad [\text{Generative Approach}] \end{array} \right.$$

The associated estimate of the pmf/pdf is therefore  $f_{\widehat{\theta}}$ , leading to the plug-in estimators

$$\left\{ \begin{array}{l} \widehat{h}^{(d)}(x) = \int_{\mathbb{R}} y f_{\widehat{\theta}}(y|x) dy \text{ for regression} \\ \widehat{h}^{(d)}(x) = \arg \max_{k=1, \dots, K} f_{\widehat{\theta}}(k|x) \text{ for classification} \end{array} \right.$$

<sup>1</sup>We use the same notation  $\widehat{R}_n(h)$  for the random variable  $\frac{1}{n} \sum_{i=1}^n L(Y_i, h(X_i))$  and its realisation  $\frac{1}{n} \sum_{i=1}^n L(y_i, h(x_i))$ .

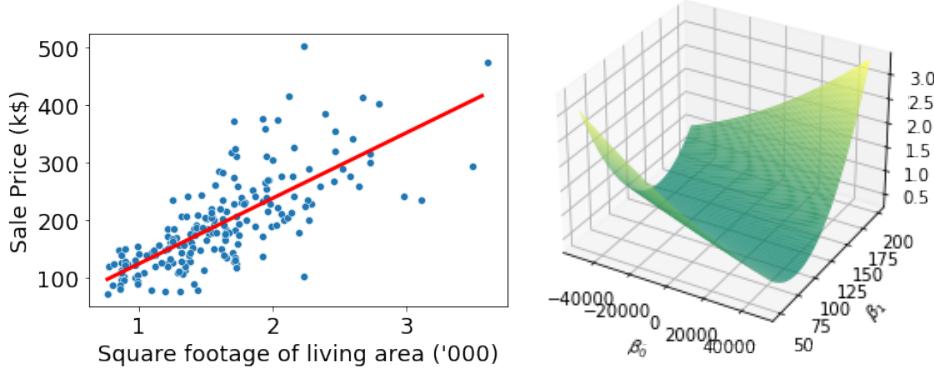


Figure 1.2: (Left) House price (dollars) versus living area (square feet) (blue dots) and estimated prediction rule  $\hat{h}^{(d)}(x)$  under a linear model (red). (Right) Empirical Risk as a function of the model parameters  $\beta_0$  and  $\beta_1$ .

### 1.4.3 Example: ordinary linear regression

As an illustration we will consider the familiar example of univariate linear regression.

We consider a dataset consisting of house sale prices together with various attributes (square footage, number of rooms, etc)<sup>2</sup>. The objective is to predict the sale price of an house based on its attributes. We will only consider here predicting the sale price  $y \in \mathbb{R}$  based on the square footage  $x \in \mathbb{R}$  of the living area. We take a training set of  $n = 200$  observations  $(x_i, y_i)$ ,  $i = 1, \dots, n$ . The data are represented in Figure 1.2.

**Empirical Risk minimisation.** Consider the set of linear prediction rules

$$\mathcal{H}_1 = \{h : \mathcal{X} \rightarrow \mathcal{Y} \mid h(x) = \beta_0 + \beta_1 x, \text{ with } \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}\}.$$

For a prediction rule  $h \in \mathcal{H}_1$ , with  $h(x) = \beta_0 + \beta_1 x$ , the empirical risk under the squared error loss is given by

$$\widehat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x)^2.$$

The empirical risk is plotted in Figure 1.2 (right) as a function of the two parameters  $\beta_0$  and  $\beta_1$ .

The ERM under the squared error loss, in the hypothesis set  $\mathcal{H}_1$ , is therefore given by

$$\hat{h}^{(d)}(x) = \hat{\beta}_0 + \hat{\beta}_1 x \tag{1.32}$$

where

$$\underline{(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(\beta_0, \beta_1) \in \mathbb{R}^2} \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x)^2}$$

are the ordinary least squares estimates, whose expression is given by Equation 1.6.

 **Plug-in approach.** In this particular case, the estimated prediction rule (1.32) may also be interpreted as a plug-in estimator, under additional assumptions. Assume now that the conditional distribution of  $Y$  given  $X = x$  is normal with mean  $\beta_0 + \beta_1 x$  and variance  $\sigma^2$ , where  $\sigma^2$  is assumed fixed. The conditional pdf is therefore given by

$$f_\theta(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\beta_0-\beta_1 x)^2}{2\sigma^2}}$$

where  $\theta = (\beta_0, \beta_1) \in \mathbb{R}^2$  are unknown parameters. The log-likelihood takes the form

$$\ell(\beta_0, \beta_1) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

<sup>2</sup><https://www.kaggle.com/c/house-prices-advanced-regression-techniques/data>

The maximum likelihood estimator is given by

$$\begin{aligned}(\hat{\beta}_0, \hat{\beta}_1) &= \arg \max_{(\beta_0, \beta_1)} -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \\&= \arg \min_{(\beta_0, \beta_1)} \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\end{aligned}$$

which is the ordinary least square estimator. The learned prediction rule is therefore

$$\widehat{h}^{(d)}(x) = \int_{\mathbb{R}} y f_{(\hat{\beta}_0, \hat{\beta}_1)}(y|x) dy = \hat{\beta}_0 + \hat{\beta}_1 x$$

which is the same as the ERM.

## 2 — Generative classifiers

$\frac{1}{2}$

This chapter is concerned with generative approaches for multiclass classification. We will present three generative classifiers: Linear Discriminant Analysis, Quadratic Discriminant Analysis and Naive Bayes. We start by a general definition of a generative classifier.

### 2.1 General definition

As mentioned in the previous chapter, generative classifiers attempt to obtain some estimate of the joint distribution of the pair  $(X, Y) \in \mathcal{X} \times \{1, \dots, K\}$ . Let

- $\pi_k = \Pr(Y = k)$  be the class probability and
- $g_k(x)$  be the conditional pdf or pmf of  $X$ , given  $Y = k$ .

An application of Bayes' theorem gives

$$\Pr(Y = k | X = x) = \frac{\pi_k g_k(x)}{\sum_{\ell=1}^K \pi_\ell g_\ell(x)}. \quad (2.1)$$

Under a  $0 - 1$  loss, the optimal classifier is

$$\begin{aligned} h^*(x) &= \arg \max_{k \in \{1, \dots, K\}} \Pr(Y = k | X = x) \\ &= \arg \max_{k \in \{1, \dots, K\}} \pi_k g_k(x) \end{aligned}$$

as the denominator in Equation (2.1) does not depend on  $k$ .

**Definition 11 (Generative classifier).** For  $k = 1, \dots, K$  and  $x \in \mathcal{X}$ , let  $\hat{\pi}_k$  and  $\hat{g}_k(x)$  be some estimates, using the training set  $d$ , of the class probabilities and of the conditional densities. The generative classifier is then given by

$$\begin{aligned} \hat{h}^{(d)}(x) &= \arg \max_{k \in \{1, \dots, K\}} \hat{\pi}_k \hat{g}_k(x) \\ &= \arg \max_{k \in \{1, \dots, K\}} \log(\hat{\pi}_k) + \log(\hat{g}_k(x)) \end{aligned}$$

$$\mathcal{L}(\pi_1, \dots, \pi_K) = \prod_{k=1}^K \pi_k^{m_k}$$

How do we obtain these estimates? For the class probabilities, it is rather straightforward. The class variables  $Y_1, \dots, Y_n$  are iid from a discrete distribution on  $\{1, \dots, K\}$  with probability mass function  $(\pi_1, \dots, \pi_K)$ . We can estimate  $\pi_k$  by maximum likelihood. The log-likelihood is given by

$$\ell(\pi_1, \dots, \pi_K) = \sum_{k=1}^K m_k \log(\pi_k) \quad (2.2)$$

where  $m_k = \#\{j : y_j = k\}$  is the number of observations in the class  $k$ . We aim at maximising the log-likelihood in (2.2) under the constraints  $\pi_k \geq 0$  for all  $k \geq 0$  and  $\sum_k \pi_k = 1$ . The Lagrangian is

$$\mathcal{L}(\pi_1, \dots, \pi_K, \gamma) = \sum_{k=1}^K m_k \log(\pi_k) - \gamma \left( \sum_{k=1}^K \pi_k - 1 \right). \quad \frac{m_k}{\pi_k} - \gamma = 0$$

Differentiating  $\mathcal{L}$  with respect to  $\pi_k$  and setting to 0 gives  $\pi_k = m_k / \gamma$  for all  $k$ . Combining this with the constraint  $\sum_{k=1}^K \pi_k = 1$ , we obtain the maximum likelihood estimate

$$\hat{\pi}_k = \frac{m_k}{n}.$$

We now describe three different methods for obtaining estimates of  $g_k(x)$ , leading to different classifiers.

$$\varphi(x; \mu_k, \Sigma) = (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right]$$

## 2.2 Linear Discriminant Analysis

Let  $\mathcal{X} = \mathbb{R}^p$ . Linear discriminant analysis (LDA) assumes the class conditional densities  $g_k(x)$  are normal pdfs, with a shared covariance matrix

$$\underbrace{g_k(x)}_{\varphi(x; \mu_k, \Sigma)}$$

where  $\mu_k$ ,  $k = 1, \dots, p$  is a vector of size  $p$ , called the centroid of the class, and  $\Sigma$  is a  $p \times p$  covariance matrix. Note that under this assumption,  $\Sigma$  is the within-class covariance matrix of the random vector  $X$ .

The parameters  $\mu_k$  and  $\Sigma$  can be estimated via maximum likelihood, as for the class frequencies  $\pi_k$ . Assuming that the  $(x_i, y_i)$  are iid realisations of a random variable  $(X, Y)$  with  $\Pr(Y = k) = \pi_k$  and conditional class densities  $g_k(x)$  given by 2.3 the log-likelihood takes the form

$$\ell(\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma) = \sum_{i=1}^n (\log(\pi_{y_i}) + \log \varphi(x_i; \mu_{y_i}, \Sigma))$$

$$= \sum_{k=1}^K \left( m_k \log(\pi_k) + \sum_{i|y_i=k} \log \varphi(x_i; \mu_k, \Sigma) \right) \quad (2.4)$$

$$= \underbrace{\left( \sum_{k=1}^K m_k \log(\pi_k) \right)}_{l_1(\pi_1, \dots, \pi_K)} - \underbrace{\frac{1}{2} \sum_{k=1}^K \sum_{i|y_i=k} (x_i - \mu_k)^\top \Sigma^{-1} (x_i - \mu_k)}_{l_2(\mu_1, \dots, \mu_K, \Sigma)} - \underbrace{\frac{n}{2} \log |\Sigma| + \text{const}}_{l_3(\Sigma)} \quad (2.5)$$

$$l_1(\pi_1, \dots, \pi_k) \quad l_2(\mu_1, \dots, \mu_k, I) \quad l_3(\Sigma) \quad (2.3)$$

Maximising the log-likelihood under the constraint  $\sum_{k=1}^K \pi_k = 1$  gives the maximum likelihood estimates

$$\left\{ \begin{array}{l} \hat{\pi}_k = \frac{m_k}{n} \\ \hat{\mu}_k = \frac{1}{m_k} \sum_{i:y_i=k} x_i \\ \hat{\Sigma} = \frac{1}{n} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^\top. \end{array} \right. \quad (2.8)$$

For a  $p$ -by- $p$  matrix  $\Sigma$ , let  $|\Sigma|$  denote the determinant of  $\Sigma$ . Let  $a, b \in \mathbb{R}^p$ . For a full-rank symmetric matrix  $\Sigma$ ,

$$\frac{\partial \log |\Sigma|}{\partial \Sigma} = \underline{\Sigma^{-1}} \quad (5)$$

$$\frac{\partial (a^\top \Sigma^{-1} b)}{\partial \Sigma} = \underline{\Sigma^{-1} a b^\top \Sigma^{-1}} \quad (6)$$

*Proof.* The MLE for  $\pi$  has already been derived. For  $\mu_k$ , the partial derivative is

$$\frac{\partial \ell}{\partial \mu_k} = \sum_{i|y_i=k} \Sigma^{-1}(x_i - \mu_k) \quad \quad \frac{\partial l}{\partial \mu_k} = 0 \quad \Rightarrow \hat{\mu}_k = \frac{1}{m_k} \sum_{i|y_i=k} x_i$$

Setting this to 0 yields the MLE (2.7). Differentiating with respect to  $\Sigma$  (see Section 0.2), we obtain

$$\frac{\partial \ell}{\partial \Sigma} = \frac{1}{2} \sum_{k=1}^K \sum_{i|y_i=k} \Sigma^{-1}(x_i - \mu_k)(x_i - \mu_k)^\top \Sigma^{-1} - \frac{n}{2} \Sigma^{-1}.$$

Setting this to 0, then left and right-multiplying by  $\Sigma$  gives Equation (2.8).

One then takes  $\hat{\pi}_k \varphi(x; \hat{\mu}_k, \hat{\Sigma})$  as an approximation to  $\pi_k g_k(x)$ . As the class densities are normal, this leads to linear discriminant functions and linear decision boundaries, as described in Section 1.3.

$$\hat{h}^{(d)}(x) = \arg \max_{k=1,\dots,K} \log(\hat{\pi}_k) - \frac{1}{2} \underbrace{(x - \hat{\mu}_k)^\top \hat{\Sigma}^{-1} (x - \hat{\mu}_k)}_{\text{squared Mahalanobis distance}}$$

$$= \arg \max_{k=1,\dots,K} \underbrace{\log(\hat{\pi}_k) - \frac{1}{2} \hat{\mu}_k^\top \hat{\Sigma}^{-1} \hat{\mu}_k}_{\hat{a}_k} + \underbrace{\hat{\mu}_k^\top \hat{\Sigma}^{-1} x}_{\hat{b}_k^\top}$$

The scalar  $\hat{a}_k$  is the intercept of the linear prediction rule for class  $k$ , and the  $p$ -dimensional vector  $\hat{b}_k$  its coefficients (called slope if  $p = 1$ ).

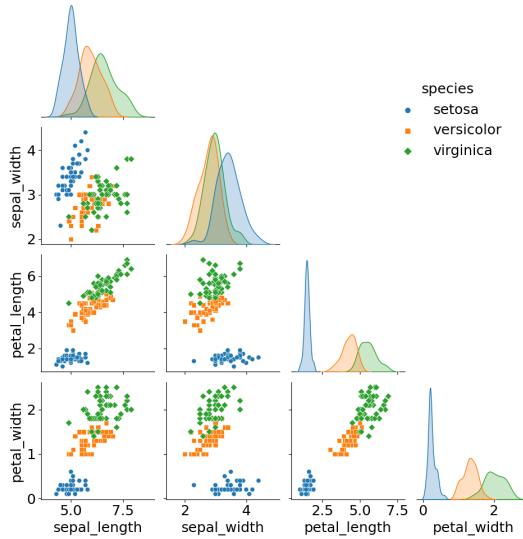


Figure 2.1: Pairplot of the Iris dataset.

**Computations for LDA.** Assuming  $\widehat{\Sigma}$  has full rank, let  $\widehat{\Sigma} = V\Lambda V^\top$  be the eigenvalue decomposition of the matrix  $\widehat{\Sigma}$ , where  $V$  is a  $p \times p$  orthogonal matrix and  $\Lambda$  is a diagonal  $p \times p$  matrix. Note that  $\widehat{\Sigma}^{-1} = V\Lambda^{-1}V^\top$  and  $\widehat{\Sigma}^{-1/2} = \Lambda^{-1/2}V^\top$ . We have

$$(x - \widehat{\mu}_k)^\top \widehat{\Sigma}^{-1} (x - \widehat{\mu}_k) = (\Lambda^{-1/2}V^\top x - \Lambda^{-1/2}V^\top \widehat{\mu}_k)^\top (\Lambda^{-1/2}V^\top x - \Lambda^{-1/2}V^\top \widehat{\mu}_k) \quad (2.9)$$

$$\sqrt{\Delta}^\top V^\top = (\Lambda^{-1/2}V^\top)^\top (\Lambda^{-1/2}V^\top) \quad (2.10)$$

that is the squared Euclidean distance between the transformed vectors  $x^* = \Lambda^{-1/2}V^\top x$  and  $\widehat{\mu}_k^* = \Lambda^{-1/2}V^\top \widehat{\mu}_k$ . Based on the estimated  $V$ ,  $\Lambda$  and  $\widehat{\mu}_k^*$ , LDA can therefore be implemented via the following steps.

- Training
  - Compute the MLE  $\widehat{\pi}_k$ ,  $\widehat{\mu}_k$  and  $\widehat{\Sigma}$ .
  - Compute the eigendecomposition  $V$ ,  $\Lambda$  of  $\widehat{\Sigma}$
  - For  $k = 1, \dots, K$ , set  $\widehat{\mu}_k^* = \Lambda^{-1/2}V^\top \widehat{\mu}_k$ .
- Prediction. For an input vector  $x$ 
  - Calculate  $x^* = \Lambda^{-1/2}V^\top x$
  - Classify to the closest class mean  $\widehat{\mu}_k^*$ , modulo the effect of the estimated class proportions  $\widehat{\pi}_k$

$$\widehat{h}^{(d)}(x) = \arg \min_{k=1, \dots, K} \|x^* - \widehat{\mu}_k^*\|^2 - 2 \log(\widehat{\pi}_k)$$

**Example 12** (Iris dataset (2D)). We consider here the classic Iris dataset. The original dataset contains  $p = 4$  measurements (sepal length, sepal width, petal length, petal width) for  $n = 50$  Iris plants. Each plant is from one of  $K = 3$  types: Iris Setosa, Iris Versicolour and Iris Virginica. A pairplot of the data is represented in Figure 12

$K-1 = 2$  discriminant coordinates.

For visualisation purposes, we will first consider two dimensions, the petal width and length. Figure 2.2 shows the scatterplot of the data, together with the LDA learned prediction rule on the two-dimensional data.

**LDA for dimensionality reduction.** Let  $\widehat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$  be the sample mean and  $\widehat{B} = \sum_{k=1}^K \widehat{\pi}_k (\widehat{\mu}_k - \widehat{\mu})(\widehat{\mu}_k - \widehat{\mu})^\top$  be an estimate of the between class covariance. Define

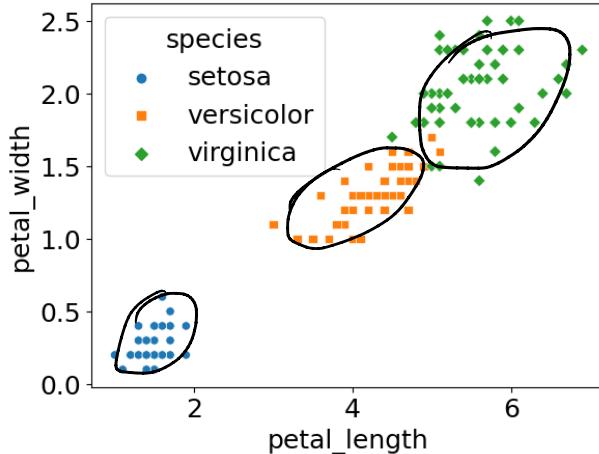
$$\widehat{B}^* = \sum_{k=1}^K \widehat{\pi}_k (\widehat{\mu}_k^* - \widehat{\mu}^*) (\widehat{\mu}_k^* - \widehat{\mu}^*)^\top$$

where  $\widehat{\mu}^* = \sum_{k=1}^K \widehat{\pi}_k \mu_k^*$ . Denote

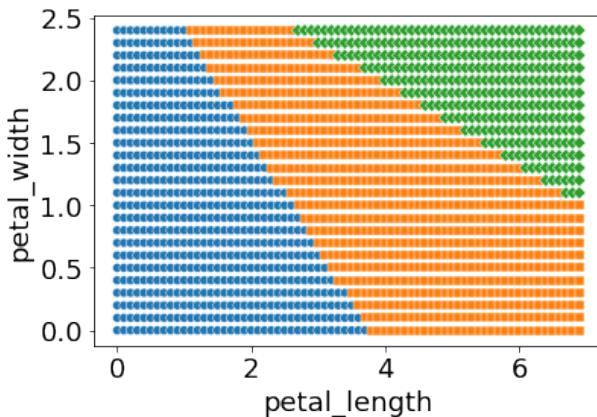
$$\widehat{B}^* = U^* D^* U^{*\top}$$

$$\widehat{\beta}^* = \Delta^{-1/2} V^\top \beta V \Delta^{-1/2}$$

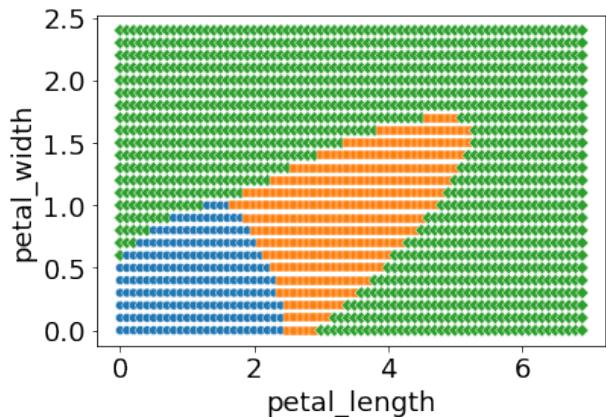
$$\widehat{\alpha}^* = \Delta^{-1/2} V^\top$$



(a) Iris data (2D) - training dataset



(b) Iris data (2D) - LDA prediction rule



(c) Iris data (2D) - QDA prediction rule

Figure 2.2: (a) Pair Scatterplot of the iris dataset in two dimensions, (b) LDA and (c) QDA prediction rules

the eigendecomposition of  $\widehat{B}^\bullet$ , and  $u_\ell^\bullet$  the columns of  $U^\bullet$ . The  $\ell$ 's **discriminant coordinate** of a vector  $x$  is defined as

$$z_\ell = (u_\ell^\bullet)^\top x^\bullet = (u_\ell^\bullet)^\top \Lambda^{-1/2} V^\top x.$$

$$\alpha_l = V \Lambda^{-\frac{1}{2}} \cdot u_l^\bullet$$

As mentioned in Section 1.3, the prediction  $\widehat{h}^{(d)}(x)$  only depends on  $x$  through its projections onto the first  $K - 1$  discriminant coordinates. The projections maximise the separation between the classes. In particular,  $a_1 = V \Lambda^{-1/2} u_1^\bullet$  minimises

$$\frac{a_1^\top \widehat{B} a_1}{a_1^\top \widehat{\Sigma} a_1}$$

$\widehat{\beta}$ : estimate of between class cov.  
 $\widehat{\Sigma}$ : estimate of within class cov.

and  $a_1^\top x$  is the one-dimensional projection of the data maximising the separation between classes.

**Example 13** (Iris dataset (4D)). We consider again the Iris dataset, with all  $p = 4$  features. As  $K = 3$ , the number of discriminant coordinates is  $K - 1 = 2$ , and we can represent the projections of the data on the plane. The projection of the Iris data on the two LDA components is represented in Figure 2.3. Note that the first discriminant coordinate already provide a good separation between the classes. For comparison, the projections onto the first two PCA components are also given. PCA maximises the variance, while LDA maximises the separation between the classes, leading to a better separation of the classes.

**LDA and SVD.** Let  $\widetilde{\mathbf{X}}$  be the  $n$ -by- $p$  matrix with rows  $(x_i - \widehat{\mu}_{y_i})^\top$ . Note that  $\widehat{\Sigma} = \frac{1}{n} \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}}$ . Similarly to the PCA case, one can use the singular value decomposition of the  $\widetilde{\mathbf{X}}$  instead of the eigendecomposition of  $\widehat{\Sigma}$  to compute the different quantities of interest. This is particularly relevant when  $p$  is large. In the case  $p > n$  the matrix  $\widehat{\Sigma}$  is not invertible anymore; we can still compute the discriminant coordinates using the SVD decomposition. We give an illustration below on an image dataset.

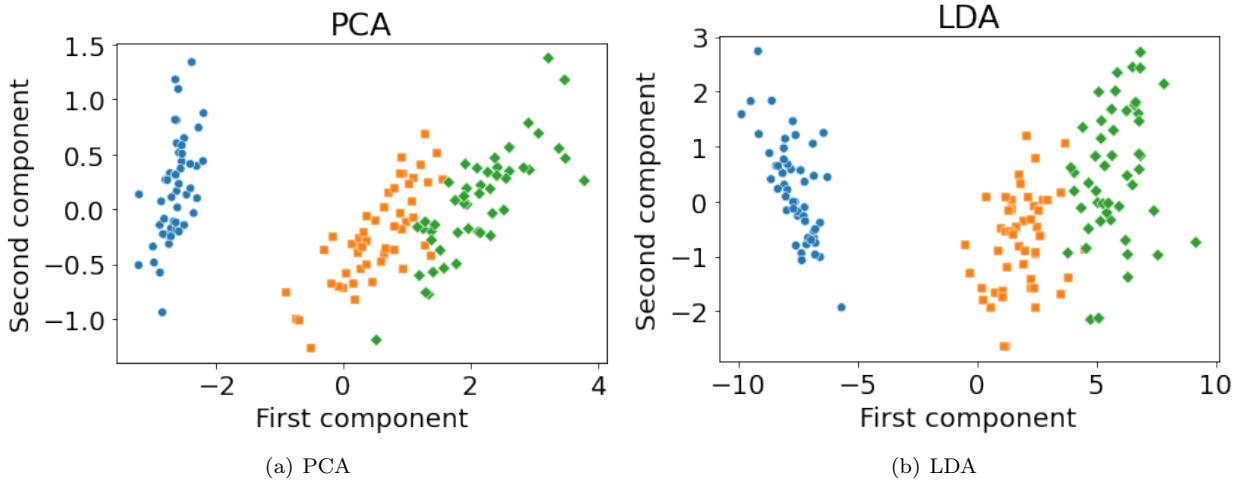


Figure 2.3: Projection of the Iris data ( $n = 150, p = 4$ ) on the first two (a) principal components of PCA and (b) discriminant coordinates of LDA.

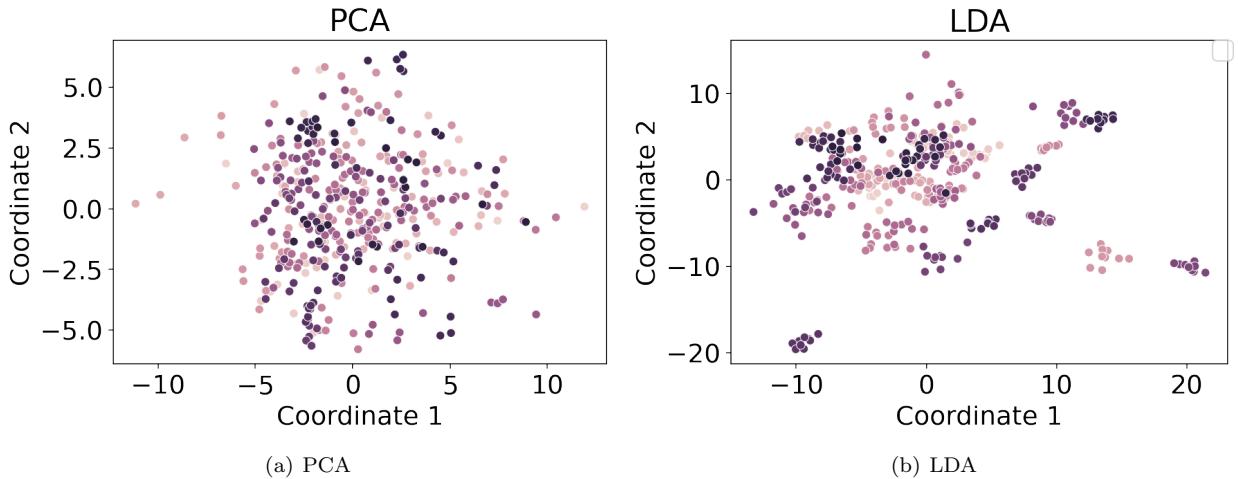


Figure 2.4: Projection of the Olivetti data ( $n = 400, p = 4096$ ) on the first two (a) principal components of PCA and (b) discriminant coordinates of LDA. The different colors represent different individuals ( $K = 40$ )

**Example 14** (Olivetti face dataset). We revisit the Olivetti dataset, that was analysed using Principal Component Analysis. The dataset contains  $n = 400$  images from  $K = 40$  different persons. Each image is  $64 \times 64$  pixel in grey levels, hence  $p = 4096$ . Note that  $p > n$  here. The first two discriminant coordinates of the training data are represented in Figure 2.4.

## 2.3 Quadratic Discriminant Analysis

Quadratic discriminant analysis (QDA) assumes the class conditional densities  $g_k(x)$  are normal pdfs, with unconstrained covariance matrices

$$g_k(x) = \underline{\varphi(x; \mu_k, \Sigma_k)}$$

where  $\mu_k, k = 1, \dots, p$  is a vector of size  $p$ , called the centroid of the class, and  $\Sigma_k$  is a  $p \times p$  covariance matrix of the class  $k$ .

The parameters  $\mu_k$  and  $\Sigma_k$  can be estimated via maximum likelihood, as for the class frequencies  $\pi_k$ . Assuming that the  $(x_i, y_i)$  are iid realisations of a random variable  $(X, Y)$  with  $\Pr(Y = k) = \pi_k$  and conditional

class densities  $g_k(x)$  given by [2.11] the log-likelihood takes the form

$$\ell(\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K) \quad (2.11)$$

$$= \sum_{i=1}^n (\log(\pi_{y_i}) + \log \varphi(x_i; \mu_{y_i}, \Sigma_{y_i})) \quad (2.12)$$

$$= \sum_{k=1}^K \left( m_k \log(\pi_k) + \sum_{i|y_i=k} \log \varphi(x_i; \mu_k, \Sigma_k) \right) \quad (2.13)$$

$$= \underbrace{\left( \sum_{k=1}^K m_k \log(\pi_k) \right)}_{l_1(\pi_1, \dots, \pi_K)} - \underbrace{\frac{1}{2} \sum_{k=1}^K \sum_{i|y_i=k} (x_i - \mu_k)^\top \Sigma_k^{-1} (x_i - \mu_k)}_{l_2(\mu_1, \dots, \mu_K, \Sigma)} - \underbrace{\sum_{k=1}^K \frac{m_k}{2} \log |\Sigma_k|}_{l_3(\Sigma_k)} + \text{const} \quad (2.14)$$

Maximising the log-likelihood under the constraint  $\sum_{k=1}^K \pi_k = 1$  gives the maximum likelihood estimates

$$\hat{\pi}_k = \frac{m_k}{n} \quad (2.15)$$

$$\hat{\mu}_k = \frac{1}{m_k} \sum_{i:y_i=k} x_i \quad (2.16)$$

$$\hat{\Sigma}_k = \frac{1}{m_k} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T. \quad (2.17)$$

One then takes  $\hat{\pi}_k \varphi(x; \hat{\mu}_k, \hat{\Sigma}_k)$  as an approximation to  $\pi_k g_k(x)$ . This leads to quadratic discriminant functions and quadratic decision boundaries, as described in Section 1.3:

$$\begin{aligned} \hat{h}^{(d)}(x) &= \arg \max_{k=1, \dots, K} \log(\hat{\pi}_k) - \frac{1}{2} \underbrace{(x - \hat{\mu}_k)^\top \hat{\Sigma}_k^{-1} (x - \hat{\mu}_k)}_{\text{squared Mahalanobis distance}} \\ &= \arg \max_{k=1, \dots, K} \underbrace{\log(\hat{\pi}_k)}_{\hat{\alpha}_k} - \underbrace{\frac{1}{2} \hat{\mu}_k^\top \hat{\Sigma}_k^{-1} \hat{\mu}_k}_{\hat{b}_k^\top x} + \underbrace{\hat{\mu}_k^\top \hat{\Sigma}_k^{-1} x}_{x^\top \hat{\Sigma}_k^{-1} x} - \underbrace{\frac{1}{2} x^\top \hat{\Sigma}_k^{-1} x}_{x^\top \hat{\alpha}_k x}. \end{aligned}$$

The QDA classifier on the 2D iris dataset is shown in Figure 2.2(c).

## 2.4 Regularised discriminant analysis

QDA allows for more flexibility and can capture more complex distributions as it allows for different covariance matrices. There is however a price to pay in terms of increased variance and potential overfitting. The number of parameters to estimate is

- $(K-1) + Kp + p(p+1)/2$  for LDA,
- $(K-1) + Kp + Kp(p+1)/2$  for QDA.

The number of parameters therefore scales quadratically with  $p$ . When  $p$  is large, this may lead to overfitting, in particular for QDA if some classes have few examples.

**Regularised Discriminant Analysis.** Regularised discriminant analysis combines LDA and QDA by considering the following estimate for the covariance within class  $k$

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}$$

where  $\hat{\Sigma}$  is the MLE for LDA, and  $\hat{\Sigma}_k$  is the MLE for QDA, for some  $\alpha \in [0, 1]$ . This introduces a new parameter  $\alpha$  and allows for a continuum of models between LDA and QDA to be used. The parameter  $\alpha$  can be chosen by cross-validation for example.

**LDA with shrinkage.** An alternative approach is to consider the shrinkage estimate (here for LDA)

$$\hat{\Sigma}(\delta) = \delta I_p + (1-\delta) \hat{\Sigma} \quad \begin{matrix} \text{can also be done} \\ \text{for QDA} \end{matrix} \quad (2.18)$$

where  $\delta \in (0, 1)$  is a regularising parameter. Note that  $\hat{\Sigma}(\delta)$  is always invertible, even if  $\hat{\Sigma}$  is not.

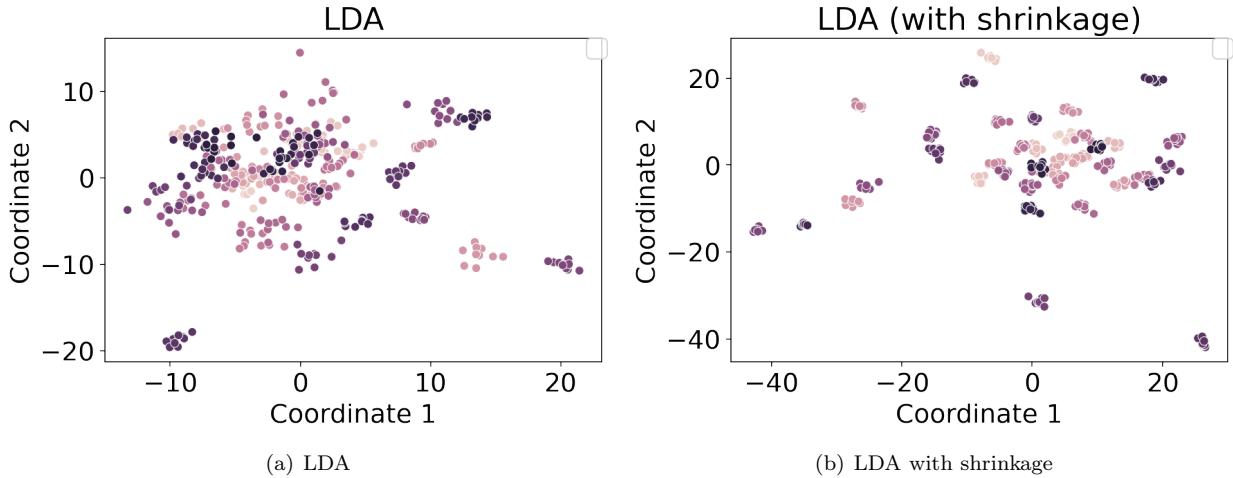


Figure 2.5: Olivetti data ( $n = 400, p = 4096$ ). First two discriminant coordinates for (left) LDA and (right) LDA with shrinkage. The different colors represent different individuals ( $K = 40$ ).

**Example 15** (Olivetti face dataset, continued). We consider again the Olivetti dataset, and look at the discriminant coordinates when using LDA, and LDA with shrinkage estimate (2.18) where the shrinkage parameter  $\delta$  is chosen automatically. The results are shown in Figure (2.5). Using shrinkage provides much improved separation between the classes for the discriminant coordinates. Note that for the LDA, coordinates are obtained using the SVD, while for LDA with shrinkage it is using the eigendecomposition, which is much slower due to the large value  $p$ .

## 2.5 Naive Bayes classifier

The naive Bayes classifier is another generative plug-in classifier. It makes the assumption that the different univariate components of  $p$ -dimensional input variable  $X$  are conditionally independent given the class  $Y$ . Naive Bayes can easily handle inputs  $x$  of mixed type (categorical, binary, continuous), as well as missing data.

To illustrate the method, consider the following example. Imagine we want to predict the voter preferences for the US election based on some attributes. Consider the data

Voted in 2016?	Annual Income	State	Candidate Choice
Y	50K	OK	Biden
N	173K	CA	Biden
Y	80K	NJ	Trump
Y	150K	WA	Biden
Y	85K	IL	Trump
:	:	:	:
Y	1050K	NY	Biden
N	35K	CA	Trump
<b>N</b>	<b>100K</b>	<b>NY</b>	<b>?</b>

The input vector  $x$  is 3-dimensional and contains a mix of binary (Voted in 2016?), real (Annual Income) and categorical (State) variables. The response variable is binary (Biden/Trump).

The naive Bayes classifier assumes that the different inputs coordinates are independent conditional on the class labels. That is, for an input vector  $\mathbf{x} = (x_1, \dots, x_p)^\top$ , the conditional class distributions, for  $k = 1, \dots, K$ , factorise as

$$g_k(x) = \prod_{j=1}^p g_{kj}(x_j; \theta_{kj}) \quad (2.19)$$

---

<sup>1</sup>We use an abuse of notation, where  $x_j$  denotes the  $j$  one-dimensional component of the vector  $x$ . Note to be confused with  $x_i$ , which will denote the  $i$ th  $p$ -dimensional observation.

where  $g_{kj}$  are conditional pmf/pdf on  $\mathbb{R}$  of  $X_j \mid Y = k$ , parameterised by  $\theta_{kj}$ . Clearly, the independence assumption is “naive” and never satisfied. But this assumption allows to significantly reduce the number of parameters to estimate and make inference much easier. Although the generative model is quite simple, naive Bayes often works quite well in practice for predicting the class labels.

In the above example, we need to estimate  $\Pr(\text{person } x \text{ voted in 2016} \mid \text{Biden supporter})$ , the probability that someone voted in 2016, given it is a Biden supporter. Similarly, we have to estimate  $\Pr(\text{person } x \text{ voted in 2016} \mid \text{Trump supporter})$ . We should do the same for the other variables “Annual income” and “State”, for each candidate, resulting in the estimation of the parameters of 6 conditional univariate pmf or pdf.

We need to define the parametric models  $g_{kj}(x_j; \theta_{kj})$ . There are standard choices depending on the type of variable.

- **Real-Valued Features.** For real-valued variables, such as the annual income, a standard choice is to use a Gaussian model, with unknown mean  $\mu_{kj}$  and variance  $\sigma_{kj}^2$ , that is

$$g_{kj}(x_j; \theta_{kj}) = \varphi(x_j; \mu_{kj}, \sigma_{kj}^2)$$

where  $\theta_{kj} = (\mu_{kj}, \sigma_{kj}^2)$ . Of course, one can use other parametric distributions.

- **Binary Features** If  $x_j \in \{0, 1\}$  (such as “Voted in 2016?”), we use a Bernoulli distribution with parameter  $\theta_{kj} \in [0, 1]$

$$g_{kj}(1; \theta_{kj}) = 1 - g_{kj}(0; \theta_{kj}) = \theta_{kj}.$$

- **Categorical Features** If  $x_j \in \{1, \dots, C\}$  is a categorical feature (such as the state), we can use the multinomial distribution, with parameters  $\theta_{kj,1}, \dots, \theta_{kj,C}$  where  $\sum_{c=1}^C \theta_{kj,c} = 1$ . For any  $c = 1, \dots, C$

$$g_{kj}(c; \theta_{kj}) = \theta_{kj,c}$$

For the special binary case  $C = 2$ , it reduces to the Bernoulli distribution.

**Parameter estimation.** We can estimate the parameters using maximum likelihood. The log-likelihood takes the form

$$\begin{aligned} \cancel{\ell}((\pi_k)_{k=1,\dots,K}, (\theta_{kj})_{k=1,\dots,K; j=1,\dots,p}) &= \sum_{i=1}^n (\log(\pi_{y_i}) + \log g_{y_i}(x_i)) \\ &= \sum_{k=1}^K m_k \log(\pi_k) + \sum_{k=1}^K \underbrace{\sum_{i|y_i=k} \log g_k(x_i)}_{\ell_k(\pi_k)} \\ &= \sum_{k=1}^K m_k \log(\pi_k) + \sum_{k=1}^K \underbrace{\sum_{i|y_i=k} \sum_{j=1}^p \log g_{kj}(x_{ij}; \theta_{kj})}_{\ell_{kj}(\theta_{kj})} \\ &= \underbrace{\sum_{k=1}^K m_k \log(\pi_k)}_{\ell_1((\pi_k))} + \underbrace{\sum_{j=1}^p \sum_{k=1}^K \underbrace{\sum_{i|y_i=k} \log g_{kj}(x_{ij}; \theta_{kj})}_{\ell_{kj}(\theta_{kj})}}_{\ell_{kj}(\theta_{kj})} \end{aligned}$$

The log-likelihood therefore factorises as  $\ell_1((\pi_k)) + \sum_{j,k} \ell_{kj}(\theta_{kj})$ , and each term of the sum can be optimised separately. For  $\pi_k$ , the MLE is given as before.

$$\widehat{\pi}_k = \frac{m_k}{n}.$$

For the other parameters, the MLE is given by

$$\cancel{\theta_{kj}} = \arg \max_{\theta_{kj}} \underbrace{\sum_{i|y_i=k} \log g_{kj}(x_{ij}; \theta_{kj})}_{\ell_{kj}(\theta_{kj})}$$

and depends on the choice of the model. In particular, we have

- Gaussian likelihood:  $\widehat{\theta}_{kj} = (\widehat{\mu}_{kj}, \widehat{\sigma}_{kj}^2)$  where

$$\widehat{\mu}_{kj} = \frac{1}{m_k} \underbrace{\sum_{i|y_i=k} x_{ij}}_{\widehat{\mu}_{kj}}, \quad \widehat{\sigma}_{kj}^2 = \frac{1}{m_k} \underbrace{\sum_{i|y_i=k} (x_{ij} - \widehat{\mu}_{kj})^2}_{\widehat{\sigma}_{kj}^2}$$

- Bernoulli likelihood:

$$\hat{\theta}_{kj} = \frac{\sum_{i|y_i=k} x_{ij}}{m_k}$$

- Multinomial likelihood: For each  $c = 1, \dots, C$ ,

$$\hat{\theta}_{k,j,c} = \frac{\sum_{i|y_i=k} \mathbb{1}_{x_{ij}=c}}{m_k}$$

Once the parameters are estimated, the conditional distribution  $\Pr(Y = k | X = x)$  is approximated by

$$\cancel{x} \quad \frac{\hat{\pi}_k \prod_{j=1}^p g_{kj}(x; \hat{\theta}_{kj})}{\sum_{k'=1}^K \hat{\pi}_{k'} \prod_{j=1}^p g_{k'j}(x; \hat{\theta}_{k'j})}$$

and the naive Bayes classifier is

$$\cancel{x} \quad \hat{h}^{(d)}(x) = \arg \max_{k \in \{1, \dots, K\}} \log(\hat{\pi}_k) + \sum_{j=1}^p \log(g_{kj}(x; \hat{\theta}_{kj})).$$

**Missing data.** Going back to our example, consider that the voter did not reveal if they had voted in 2016:

Voted in 2016?	Annual Income	State	Candidate Choice
Y	50K	OK	Biden
N	173K	CA	Biden
Y	80K	NJ	Trump
Y	150K	WA	Biden
Y	85K	IL	Trump
:	:	:	:
Y	1050K	NY	Biden
N	35K	CA	Trump
?	<b>100K</b>	<b>NY</b>	?

For our new observation  $\tilde{x} = (\tilde{x}_1, \tilde{x}_2, \tilde{x}_3)$ , the value of the first variable  $\tilde{x}_1$  is missing, and we can thus only use  $(\tilde{x}_2, \tilde{x}_3)$  to predict its candidate choice. Under the naive Bayes assumption,

$$g_k((\tilde{x}_1, \tilde{x}_2, \tilde{x}_3)) = \prod_{j=1}^3 g_{kj}(\tilde{x}_j; \theta_{kj})$$

Marginalising the above expression with respect to  $\tilde{x}_1$  gives the conditional distribution of  $(\tilde{X}_2, \tilde{X}_3)$  given  $\tilde{Y} = k$ , which is

$$\cancel{O} \quad \underbrace{\sum_{\tilde{x}_1=0,1} g_k((\tilde{x}_1, \tilde{x}_2, \tilde{x}_3))}_{\text{Marginalise } \tilde{x}_1} = \prod_{j=2}^3 g_{kj}(\tilde{x}_j; \theta_{kj}). \quad (2.20)$$

Hence, using Bayes rule,

$$\Pr(\tilde{Y} = k | (\tilde{X}_2, \tilde{X}_3) = (\tilde{x}_2, \tilde{x}_3)) \propto \underbrace{\pi_k}_{\text{Prior}} \prod_{j=2}^3 g_{kj}(\tilde{x}_j; \theta_{kj}).$$

The naive Bayes prediction is therefore given by

$$\cancel{O} \quad \hat{h}^{(d)}((\tilde{x}_2, \tilde{x}_3)) = \arg \max_{k \in \{1, \dots, K\}} \log(\hat{\pi}_k) + \sum_{j=2}^p \log(g_{kj}(x; \hat{\theta}_{kj})).$$

Note that missing data can be handled in a similar way for any generative approach, by marginalising the variable of interest in  $g_k(x)$ . This is done without extra computation for naive Bayes (see Equation (2.20)) due to the independence assumption. Without the independence this would require summation/integration.

Dealing with missing data in the training set is also quite easy. Assume that some entries are missing for the variable Voted in 2016:

Voted in 2016?	Annual Income	State	Candidate Choice
?	50K	OK	Biden
?	173K	CA	Biden
Y	80K	NJ	Trump
Y	150K	WA	Biden
?	85K	IL	Trump
:	:	:	:
Y	1050K	NY	Biden
N	35K	CA	Trump
?	<b>100K</b>	<b>NY</b>	?

For example, let's say that for Biden voters, 103 had voted in 2016, 54 had not, and 25 didn't answer. We can simply estimate  $\theta_{Biden,1}$  based on the non-missing answers, that is  $\hat{\theta}_{Biden,1} = 103/157$ .

**Naive Bayes and Curse of dimensionality.** By making the independence assumption, naive Bayes breaks the curse of dimensionality and avoids overfitting. Assume for illustration that all the features  $x_j$  of the vector  $x$  are binary, hence  $x \in \{0, 1\}^p$ , where  $\{0, 1\}^p$  has  $2^p$  elements. If one were to parameterise  $g_k(x)$ , this would require  $O(2^p)$  parameters. By contrast, using naive Bayes only requires  $O(p)$  parameters (one for each dimension of the input vector  $x$ ).

**Example 16.** As an example consider the dataset consisting of attributes of 887 passengers of the Titanic, and whether or not they survived the tragedy.

Class	Sex	Age	Survived
3	male	22	N
1	female	38	Y
3	female	26	Y
1	female	35	Y
3	male	35	N
:	:	:	:

The response variable  $y$  is binary ( $\text{survived}=1$ , did not survive=0), and the input vector  $x$  is three dimensional  $x = (x_1, x_2, x_3)$  where  $x_1$  is the class of the passenger (1st, 2nd or 3rd),  $x_2$  is the sex (female/male),  $x_3$  is the age. We aim at deriving a classifier to predict the survival based on the three attributes. Histograms of the input variables for passengers who survived and did not survive are shown in Figure 2.6.

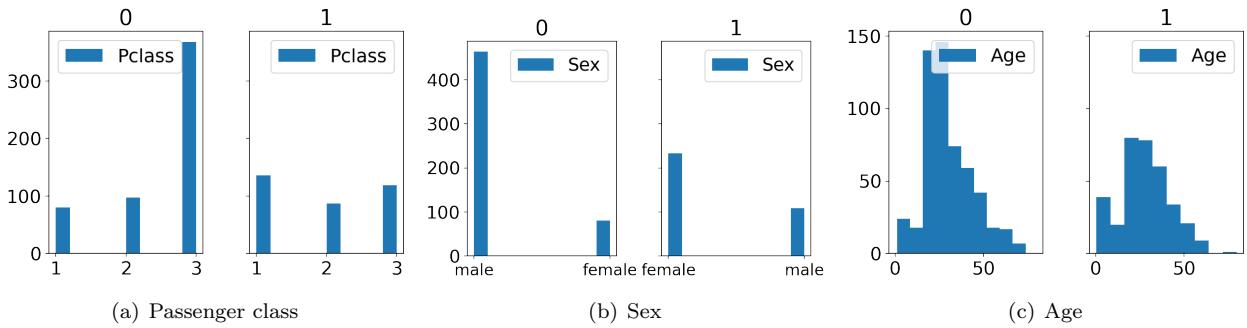


Figure 2.6: Histogram of the (a) Passenger class (b) Sex and (c) Age of the passenger amongst those who did not survive (left plot) and did survive (right plot).

We use a multinomial model with  $C = 3$  for  $g_{k1}$  (passenger class), a Bernoulli model for  $g_{k2}$  (sex) and a Gaussian model for  $g_{k3}$  (age). As can be seen from Figure 2.6(c) the Gaussian distribution provides a poor fit to the data, and one could consider alternative models. The classifier is defined by

$$\hat{h}^{(d)}(x) = \arg \max_{k \in \{0,1\}} \log(\hat{\pi}_k) + \sum_{j=1}^3 \log(g_{kj}(x; \hat{\theta}_{kj})).$$

The estimated parameters are reported in the following table

Name	Parameters	Did not Survived ( $k=0$ )	Survived ( $k=1$ )
Prior class proba	$\pi_k$	0.61	0.39
Pass. Class	$(\hat{\theta}_{k1,1}, \dots, \hat{\theta}_{k1,3})$	(0.15, 0.18, 0.67)	(0.40, 0.25, 0.35)
Sex ( $F=1$ )	$\hat{\theta}_{k2}$	0.15	0.32
Age	$(\hat{\mu}_{k3}, \hat{\sigma}_{k3})$	(30.1, 13.9)	(28.4, 14.4)

Here are some predictions from the naive Bayes classifier.

Class	Sex	Age	Predicted Survival	Estimated Probability of survival
3	female	30	Y	0.587
1	female	30	Y	0.882
3	male	20	N	0.114
1	male	20	N	0.403

## 2.6 Summary

In this chapter, we have seen three different classes of generative classifiers:

- LDA and QDA assume that the input  $X$  is normally distributed given the class  $Y = k$ , with the same covariance (LDA) or different covariances (QDA)
- Naive Bayes assumes that the different dimensions of the input  $X$  are conditionally independent given the class  $Y = k$

The advantages of using a generative approach are

- The assumptions made on the data generating process can be statistically tested (Gaussianity, independence);
- Even if the assumptions do not hold, the classifier may still provide good predictions in practice;
- The classifier can easily handle missing data in the input vector;
- It leads to interpretable predictions.

The main drawback is that it makes strong assumptions on the data generating process, that are rarely met in practice. This limits the flexibility of the generative approach.

# 3 — Further concepts in statistical learning

1/2

In this chapter, we will discuss key concepts in statistical machine learning:

- Nonlinear input transformation/expansion,
- Overfitting and bias-variance tradeoff
- Regularisation
- Cross-validation

Throughout this chapter, we will use a simple regression example to illustrate these concepts, and frame the learning procedure as an empirical risk minimisation problem over some class of prediction rules. The concepts and tools discussed in this chapter however apply equally when dealing with classification tasks, or when using plug-in methods. We start by recalling some background material on multivariate linear regression.

## 3.1 Multivariate linear regression

$$\mathbf{x} = \begin{pmatrix} \cdot \\ \vdots \\ \cdot \end{pmatrix}$$

Let  $((x_1, y_1), \dots, (x_n, y_n))$  be the input/target data, with  $x_i \in \mathbb{R}^p$ , and let  $y_i \in \mathbb{R}$ . For an input variable  $x \in \mathbb{R}^p$ , let  $\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix} \in \mathbb{R}^{p+1}$ . Consider the class of linear prediction rules

$$\mathcal{H} = \{h(x) = \tilde{x}^\top \beta = (1 \ x^\top) \beta \mid \beta \in \mathbb{R}^{p+1}\}.$$

The empirical risk minimiser under a squared error loss for the class  $\mathcal{H}$  is  $\hat{h}^{(d)}(x) = \tilde{x}^\top \hat{\beta}$  where

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (y_i - \tilde{x}_i^\top \beta)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^{p+1}} \| \mathbf{y} - \tilde{\mathbf{X}} \beta \|^2 \end{aligned}$$

where  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$  and  $\tilde{\mathbf{X}}$  is the  $n$ -by- $(p+1)$  design matrix defined by

$$\underbrace{\frac{\partial(\mathbf{x}^\top \alpha)}{\partial \mathbf{x}}} = \alpha \quad \underbrace{\frac{\partial((\mathbf{x}^\top \alpha)^\top \beta (\mathbf{x}^\top \alpha))}{\partial \mathbf{x}}} = 2\beta(\mathbf{x}^\top \alpha) \quad \tilde{\mathbf{X}} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \cdot = \begin{pmatrix} 1 & \mathbf{x}_1^\top \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top \end{pmatrix} \quad (3.1)$$

We have

$$\begin{aligned} \|\mathbf{y} - \tilde{\mathbf{X}} \beta\|^2 &= (\mathbf{y} - \tilde{\mathbf{X}} \beta)^\top (\mathbf{y} - \tilde{\mathbf{X}} \beta) \\ &= \beta^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \beta - 2(\tilde{\mathbf{X}}^\top \mathbf{y})^\top \beta + \text{const.} \end{aligned}$$

Using Equations (3) and (4), we obtain the derivative

$$\frac{\partial(\|\mathbf{y} - \tilde{\mathbf{X}} \beta\|^2)}{\partial \beta} = 2\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \beta - 2\tilde{\mathbf{X}}^\top \mathbf{y}$$

Setting this to 0, and assuming that  $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$  has full rank, we obtain the estimate

$$\hat{\beta} = \underbrace{(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbf{y}}.$$

### 3.2 Nonlinear input expansion

The learned prediction rule  $\hat{h}^{(d)}(x) = \tilde{x}^\top \hat{\beta}$  has the interesting property that it can be obtained in closed-form, but it can only capture linear relations between the input  $x$  and the response  $y$ . It is actually possible to obtain non-linear prediction rules by using a class of linear prediction rules on an extended input/feature space.

For an input space  $\mathcal{X}$ , consider a function  $\phi : \mathcal{X} \rightarrow \mathcal{X}'$ , where  $\mathcal{X}'$  is the extended input (or feature) space. The function  $\phi$  is a known, typically nonlinear function. The extended input/feature space  $\mathcal{X}'$  may be of dimension lower, equal or larger than the original input space  $\mathcal{X}$ . The function  $\phi$  may be used to reduce dimensionality and extract interesting features of the data (for instance using PCA), or to expand the dimension of the feature space to allow for more complex prediction rules. One can apply a given learning method to the transformed inputs  $\phi(x_1), \dots, \phi(x_n)$  instead of the original inputs. This potentially allows to capture more complex prediction rule under simple learning methods.

Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}^{M+1}$  for some  $M \geq 0$ . One can for example consider the class of linear prediction rules on the transformed inputs, that is

$$\mathcal{H} = \{h(x) = \phi(x)^\top \beta \text{ where } \beta \in \mathbb{R}^{M+1}\}.$$

Note that the class of prediction rules are linear functions in the transformed input space  $\mathcal{X}'$ , but nonlinear function in the original input space  $\mathcal{X}$  if the transformation  $\phi$  is nonlinear. The ERM under a squared loss can be found in closed form, as it corresponds to the ordinary least square estimate on the transformed inputs  $\phi(x_1), \dots, \phi(x_n)$ :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{M+1}} \frac{1}{n} \sum_{i=1}^n (y_i - \phi(x_i)^\top \beta)^2 = \underbrace{(\Phi^\top \Phi)^{-1} \Phi^\top}_{\Phi = \begin{pmatrix} \phi(x_1)^\top \\ \vdots \\ \phi(x_n)^\top \end{pmatrix}} \mathbf{y}$$

$$\Phi = \begin{pmatrix} \phi(x_1)^\top \\ \vdots \\ \phi(x_n)^\top \end{pmatrix}$$

where

$$\Phi = (\phi(x_1), \dots, \phi(x_n))^\top \in \mathbb{R}^{n \times (M+1)}, \quad \mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n. \quad (3.2)$$

**Example 17. (sinusoid)** As an illustrative example, consider that  $(X, Y)$  has a joint distribution defined by

$$Y = \sin(2\pi X) + \epsilon, \quad X \sim U(0, 1), \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

where  $X$  is independent of  $\epsilon$ , and  $\sigma > 0$ . The Bayes prediction rule is

$$h^*(x) = \mathbb{E}[Y | X = x] = \sin(2\pi x).$$

Excess risk of  $h$  (Risk) -  $R(h^*)$

The risk for a prediction rule  $h$  is given by

$$R(h) = \mathbb{E}[(Y - h(X))^2] = \mathbb{E}[(\sin(2\pi X) + \epsilon - h(X))^2] = \sigma^2 + \int_0^1 (\sin(2\pi x) - h(x))^2 dx$$

and the Bayes risk is  $R(h^*) = \sigma^2$ . For  $M \geq 0$ , let  $\mathcal{H}_M$  be the set of polynomials of order  $M$

$$\mathcal{H}_M = \{h : \mathcal{X} \rightarrow \mathcal{Y} \mid h(x) = \sum_{j=0}^M \beta_j x^j, \beta_j \in \mathbb{R} \text{ for } j = 0, \dots, M\}.$$

Note that the hypothesis classes are nested  $\mathcal{H}_0 \subset \mathcal{H}_1 \subset \dots$ . Note also that  $h^* \notin \mathcal{H}_M$ , for all  $M$ .

We generate  $n = 10$  realisations  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , and calculate the ERM prediction rule under the classes  $\mathcal{H}_M$ , for  $M = 0, \dots, 9$ . The learned prediction rules  $\hat{h}^{(d)}$  are plotted in Figure 3.1. For  $M = 0$  and  $M = 1$ , the learned prediction rules are too simple. We say that we are underfitting. For  $M = 3$ , the learned prediction rule is close to the Bayes prediction rule. For  $M = 9$ , the learned prediction interpolates the data, but is very far from the Bayes prediction rule when  $x$  is not in the training set. As shown in Table 3.1, the estimated parameters take very large values for  $M = 9$  suggesting the fact that the estimator has very large variance. We say that we are overfitting: the model is too complex for the amount of data available, and is fitting both the signal and the noise. Figure 3.2 shows the (generalisation/population) risk and the empirical risk of the learned prediction rule as a function of the order of the polynomial. The empirical risk decreases with the model complexity, to reach zero for  $M = 9$ . The true risk has a U shape: it first decreases (underfitting regime), then increases (overfitting regime).

Overfitting can have disastrous effect. When fitting a polynomial of degree 11 on the house pricing data, we obtain a prediction rule that decreases for large living area, as shown in Figure 3.3.

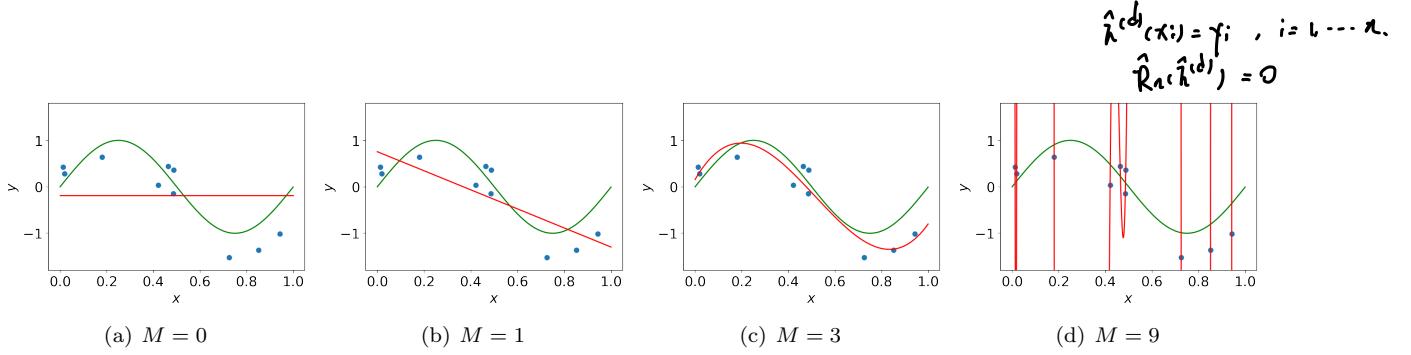
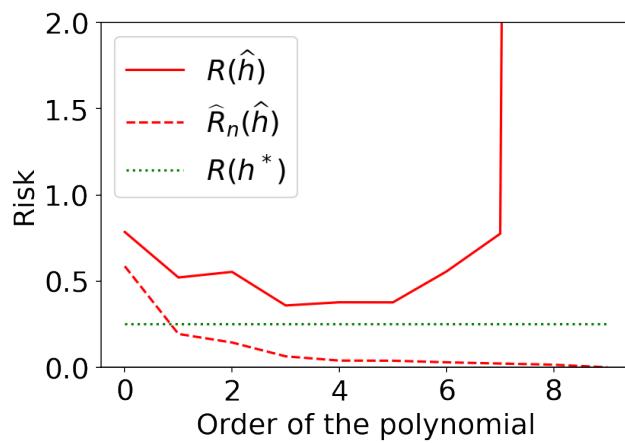


Figure 3.1: ERM prediction rules (in red) for linear models with polynomial expansion with degrees  $M = 0, 1, 3$  and  $9$ . The data are represented by blue dots and the Bayes prediction rule is in green.

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$\hat{\beta}_0$	-0.19	0.76	0.16	171.61
$\hat{\beta}_1$		-2.05	8.66	-24859.77
$\hat{\beta}_2$			-27.29	1052225
$\hat{\beta}_3$			17.67	-13060914
$\hat{\beta}_4$				75449695
$\hat{\beta}_5$				-239878655
$\hat{\beta}_6$				446010421
$\hat{\beta}_7$				-483353259
$\hat{\beta}_8$				282706993
$\hat{\beta}_9$				-68922625

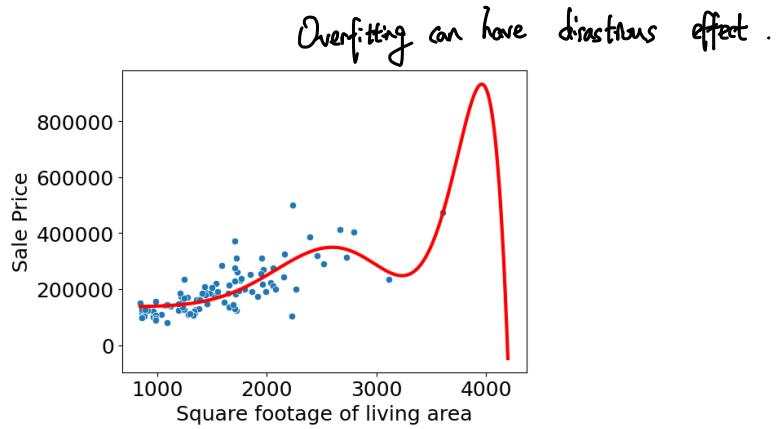
Table 3.1: Estimates under for polynomial class of degree 0, 1, 3, and 9.

Empirical risk / Training error :  $\hat{R}_n(\hat{h}) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{h}(x_i))$



$$R(\hat{h}) = \mathbb{E}[L(Y, \hat{h}(X))]$$

Figure 3.2: Generalisation Risk (plain red line) and Empirical risk (dashed red line) as a function of the model complexity. The Bayes risk is represented by a dotted green line.



$R(h^*)$ , Bayes risk

### 3.3 Estimation-Approximation and Bias-variance trade-off

**Estimation-Approximation error.** For a given set of prediction rules  $\mathcal{H} \subset \mathcal{F}$ , called an **hypothesis class**, denote

$$\underline{h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} R(h)}$$

the best prediction rule, in terms of risk minimisation, amongst the set  $\mathcal{H}$ . Note that by definition, we have  $R(\hat{h}^{(d)}) \geq R(h_{\mathcal{H}}^*) \geq R(h^*)$ , where  $\hat{h}^{(d)}$  is the empirical risk minimiser (1.30). The excess risk can be decomposed as a sum of two terms

$$\underbrace{R(\hat{h}^{(d)}) - R(h^*)}_{\text{excess risk}} = \underbrace{R(\hat{h}^{(d)}) - R(h_{\mathcal{H}}^*)}_{\text{estimation error}} + \underbrace{R(h_{\mathcal{H}}^*) - R(h^*)}_{\text{approximation error}}$$

- The **approximation error** is the difference between the risk of the best prediction rule within the class and the Bayes risk. If the hypothesis class is large, more complex prediction rules can be captured, and the approximation error is thus smaller.
- The **estimation error** is the difference between the risk of the learned prediction rule and of the best prediction rule in the class. Larger classes have a large number of parameters to estimate, leading to a larger estimation error.

This approximation-estimation trade-off and how it relates to the model complexity is illustrated in Figure 3.4

**Example 18.** (sinusoid, continued) For the simulated example described in the previous section, the model complexity corresponds to the degree of the polynomial. Figure 3.5 shows the learned prediction rules and the best prediction rule for each class. As the degree increases, the best prediction rule in the class becomes closer and closer to the Bayes prediction rule. Figure 3.6 shows the excess risk and the estimation/approximation error as a function of the degree of the polynomial.

**Approximation error and size of the dataset.** The approximation error is a quantity that only depends on the choice of the class  $\mathcal{H}$ , not the dataset. Under mild assumptions, in general we would have that the estimation error goes to 0 as the number of observations  $n$  goes to infinity. Hence the excess risk converges to the approximation error or, similarly, the risk  $R(\hat{h})$  converges to the risk  $R(h_{\mathcal{H}}^*)$  of the best prediction rule in the class as  $n \rightarrow \infty$ .

**The empirical risk underestimates the true risk.** The empirical risk of the learned rule  $\hat{R}_n(\hat{h})$  typically underestimates the true risk  $R(\hat{h})$  of that prediction rule. Note that, by definition of the ERM in the class  $\mathcal{H}$ , we have  $R(\hat{h}^{(d)}) \geq R(h_{\mathcal{H}}^*)$  for any dataset  $d$ . But (see Problem Sheet)<sup>1</sup>

$$\mathbb{E}[\hat{R}_n^{(D)}(\hat{h}^{(D)})] \leq R(h_{\mathcal{H}}^*)$$

where the expectation is taken with respect to the random sample  $D$ . The **generalisation gap** is

$$\underline{R(\hat{h}) - \hat{R}_n(\hat{h})}.$$

This generalisation gap typically converges to 0 as the sample size  $n$  increases, as both  $R(\hat{h})$  and  $\hat{R}_n(\hat{h})$  will converge to  $R(h_{\mathcal{H}}^*)$ . The relation between the generalisation gap and the sample size is depicted in Figure 3.3.

<sup>1</sup>we add the superscript  $D$  to emphasize the dependence on the random sample  $D$ .

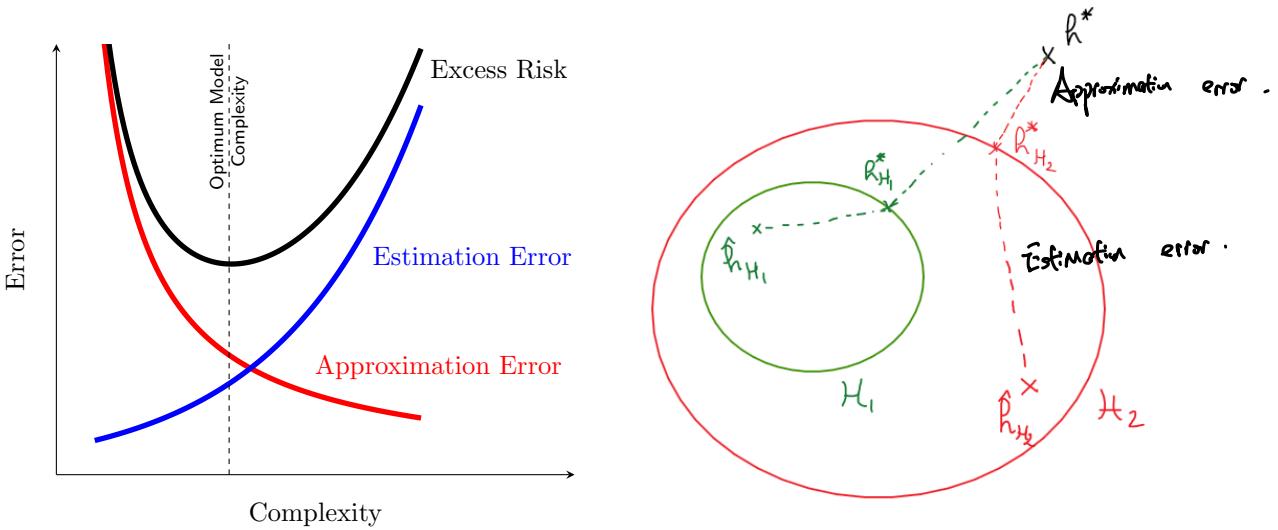


Figure 3.4: Illustration of the approximation-estimation tradeoff. (Left) Approximation and estimation error as a function of the model complexity. Simple models have large approximation error and small estimation error. Complex models have large estimation error and small approximation error. (Right) For two hypothesis classes  $\mathcal{H}_1 \subset \mathcal{H}_2$  ( $\mathcal{H}_1$  is simpler than  $\mathcal{H}_2$ ), and a Bayes rule  $h^*$  outside of the classes.  $\mathcal{H}_1$  will have larger approximation error  $R(h_{\mathcal{H}_2}^*) \leq R(h_{\mathcal{H}_1}^*)$  but typically smaller estimation error.

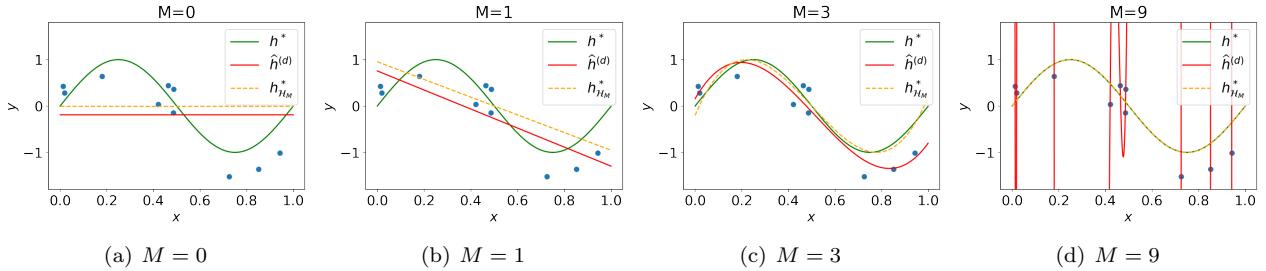


Figure 3.5: ERM learned prediction rules (in red) for linear models with polynomial expansion with degrees  $M = 0, 1, 3$  and  $9$ . The best prediction rule within the class of polynomials of order  $M$  is shown in orange. The data are represented by blue dots and the Bayes prediction rule is in green.

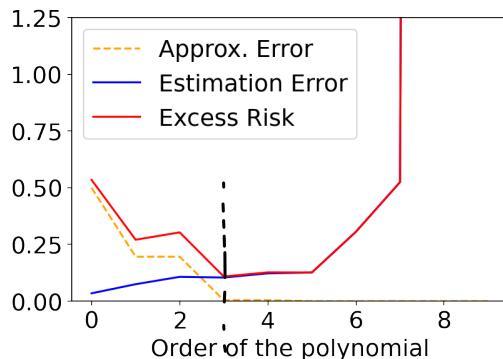


Figure 3.6: Excess risk  $R(\hat{h}^{(d)}) - R(h^*)$  (red), approximation error (orange) and estimation error (blue) as a function of the complexity.

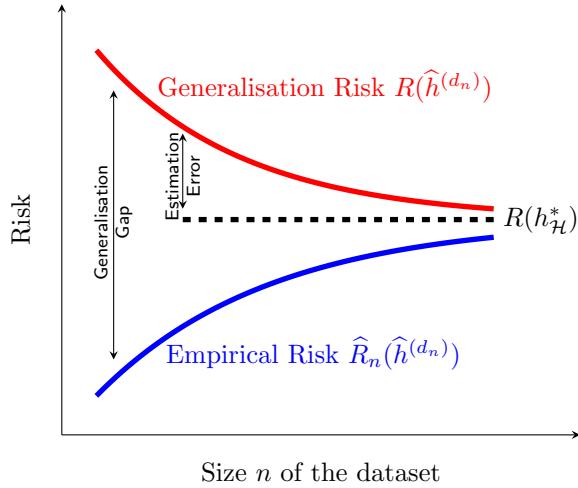


Figure 3.7: Generalisation and Empirical risk as a function of the size  $n$  of the dataset for a given hypothesis class  $\mathcal{H}$ .

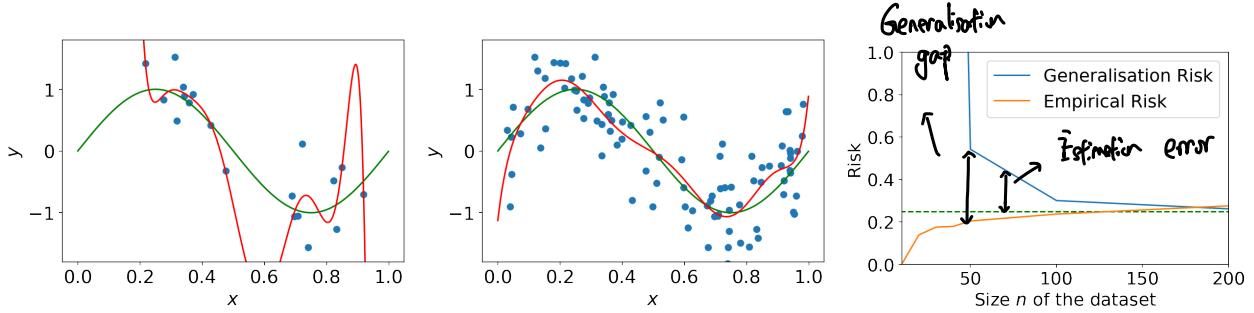


Figure 3.8: Learned prediction rule with a polynomial of order 9 with (left)  $n = 20$  and (middle)  $n = 100$  observations. (Right) Generalisation and Empirical Risks as a function of the size  $n$  of the dataset.

**Example 19.** (*sinusoid, continued*) Going back to our simulated example, consider the ERM with hypothesis class the set of polynomials of order 9, with a dataset of growing size. Figure 3.8 shows the estimated prediction rules for  $n = 20, 100$ , and the empirical and generalisation risks as a function of the sample size. Both converge to the same value, corresponding to the risk of the best prediction rule within the class of polynomials of order 9.

**Bias-variance decomposition.** For the squared loss, another traditional interpretable decomposition of the excess risk is the bias-variance decomposition. Recall that for the squared loss, the excess risk takes the form

$$\mathbb{E}_D [R(\hat{h}^{(D)}) - R(h^*)] = \mathbb{E}_D \mathbb{E}_X [(\hat{h}^{(D)}(X) - h^*(X))^2].$$

The above excess risk is a random variable, as it depends on the dataset  $D$ . Taking the expectation over the dataset  $D$ , we obtain the expected excess risk

$$\begin{aligned} \mathbb{E}_D [R(\hat{h}^{(D)}) - R(h^*)] &= \mathbb{E}_D \mathbb{E}_X [(\hat{h}^{(D)}(X) - h^*(X))^2] \\ &= \mathbb{E}_X \mathbb{E}_D [(\hat{h}^{(D)}(X) - h^*(X))^2] \end{aligned}$$

For any  $x \in \mathcal{X}$ , the random variable  $\hat{h}^{(D)}(x)$  is an estimator of the Bayes predictor  $h^*(x)$ . This estimator has mean  $\bar{h}_{\mathcal{H},n}(x) = \mathbb{E}_D [\hat{h}^{(D)}(x)]$  and bias and variance

$$\begin{cases} b_{\mathcal{H},n}(x) = \bar{h}_{\mathcal{H},n}(x) - h^*(x) \\ v_{\mathcal{H},n}(x) = \mathbb{E}_D [(\hat{h}^{(D)}(x) - \bar{h}_{\mathcal{H},n}(x))^2] \end{cases} \quad (3.3)$$

For any  $x$ , we have

$$\begin{aligned} \mathbb{E}_D [(\hat{h}^{(D)}(x) - h^*(x))^2] &= \mathbb{E}_D [(\hat{h}^{(D)}(x) - \bar{h}_{\mathcal{H},n}(x))^2 + (\bar{h}_{\mathcal{H},n}(x) - h^*(x))^2 + 2(\hat{h}^{(D)}(x) - \bar{h}_{\mathcal{H},n}(x))(\bar{h}_{\mathcal{H},n}(x) - h^*(x))] \\ &\quad (3.5) \end{aligned}$$

$$= v_{\mathcal{H},n}(x) + b_{\mathcal{H},n}(x)^2. \quad (3.6)$$

For Bias Variance decomposition:

$$R(h) = \mathbb{E} [L(Y, h(x))] , \quad L(Y, h(x)) = (Y - h(x))^2$$

$$\begin{aligned} \textcircled{1} \quad h^* &= \underset{h(x)}{\operatorname{argmin}} \mathbb{E} [(Y - h(x))^2] \\ &= \underset{h(x)}{\operatorname{argmin}} \int \mathbb{E}_x (Y - h(x))^2 | X=x f_X(x) dx \\ &= \underset{h(x)}{\operatorname{argmin}} \mathbb{E} [(Y - h(x))^2 | X=x] \\ \mathbb{E} [(Y - h(x))^2 | X=x] &= \mathbb{E} (Y^2 | X=x) - 2\mathbb{E}(Yh(x) | X=x) + \mathbb{E}(h^2(x) | X=x) \\ &= \mathbb{E} (Y^2 | X=x) - 2h \mathbb{E}(Y | X=x) + h^2 \\ &= \operatorname{Var}(Y | X=x) + [\mathbb{E}(Y | X=x) - h]^2 \\ \Rightarrow h^*(x) &= \underline{\mathbb{E}(Y | X=x)} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad R(h^*) &= \mathbb{E} [ (Y - \mathbb{E}(Y | X=x))^2 ] \\ &= \mathbb{E} [ Y^2 - 2Y\mathbb{E}(Y | X=x) + \mathbb{E}^2(Y | X=x) ] \\ &= \mathbb{E} [ \mathbb{E}^2(Y^2 | X=x) - 2\mathbb{E}(Y | X=x)\mathbb{E}(Y | X=x) + \mathbb{E}^2(Y | X=x) ] \\ &= \operatorname{Var}(Y | X) \end{aligned}$$

$$\textcircled{3} \quad R(\hat{h}) - R(h^*) = \mathbb{E} [-2\hat{h} + \hat{h}^2 + 2h^* - h^{*2}]$$

$$\mathbb{E}(Y\hat{h}) = \mathbb{E}(\mathbb{E}(Y | X)\cdot \hat{h}) = \mathbb{E}(h^*\hat{h})$$

**Solution:** We have

$$\begin{aligned} R(h) &= \mathbb{E} [(Y - h(X))^2] \\ &= \int_X \mathbb{E} [(Y - h(X))^2 | X=x] f_X(x) dx, \end{aligned}$$

where  $f_X$  is pdf/pmf of  $X$ . Thus, it suffices to minimise, for every  $x$ ,

$$\begin{aligned} &\mathbb{E} [(Y - h(X))^2 | X=x] \\ &= \mathbb{E} [Y^2 | X=x] - 2h(x)\mathbb{E}[Y | X=x] + h(x)^2 \\ &= \operatorname{var}[Y | X=x] + (\mathbb{E}[Y | X=x] - h(x))^2. \end{aligned}$$

This is clearly minimized by the conditional mean:

$$h^*(x) = \mathbb{E}[Y | X=x].$$

The associated risk is

$$R(h^*) = \mathbb{E}[\operatorname{var}(Y | X)] = \int_X \operatorname{var}(Y | X=x) f_X(x) dx$$

We have

$$\begin{aligned} R(h) &= \mathbb{E} [\operatorname{var}(Y | X) + (\mathbb{E}[Y | X] - h(X))^2] \\ &= R(h^*) + \mathbb{E}[(h(X) - h^*(X))^2]. \end{aligned}$$

$$\begin{aligned} &= \mathbb{E} (-2h^*\hat{h} + \hat{h}^2 + h^{*2}) \\ &= \mathbb{E} [(\hat{h} - h^*)^2] \end{aligned}$$

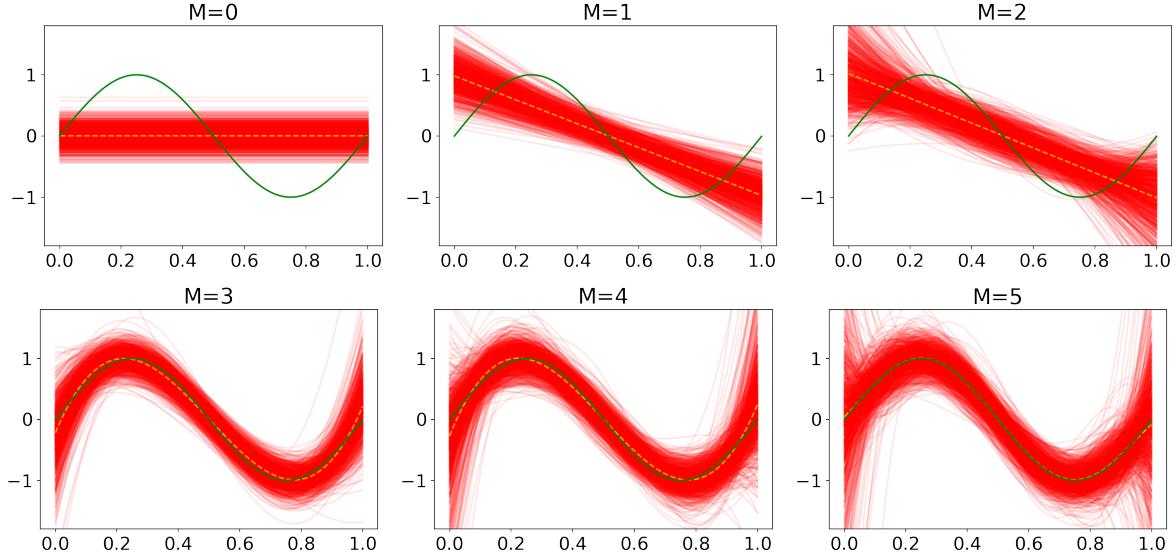


Figure 3.9: Learned ERM rules for different datasets of size  $n = 30$ , for different orders of the polynomial. As the order  $M$  increases, variance of the predictions increases, in particular near the boundary.

$$\mathbb{E}_D [\hat{h}^{(D)} - h^*] = v_{\mathcal{H},n}(x) + b_{\mathcal{H},n}(x)^2$$

It follows that

$$\mathbb{E}_D [R(\hat{h}^{(D)}) - R(h^*)] = \underbrace{\mathbb{E}_X [v_{\mathcal{H},n}(X)]}_{\text{average excess risk}} + \underbrace{\mathbb{E}_X [b_{\mathcal{H},n}(X)^2]}_{\text{variance term}} \quad (3.7)$$

- The variance term relates to the estimation error; it will typically be large for large hypothesis sets  $\mathcal{H}$ . This term vanishes as the number of observations grows.
- The bias term relates to the approximation error; it will typically be large for small hypothesis sets  $\mathcal{H}$ . This term does not vanish as the number of observations grows.

**Example 20.** (sinusoid, continued) Using the same simulated example, with  $n = 30$  examples, we sample different datasets and represent for each dataset the learned prediction rule in Figure 3.3. The mean of these learned prediction rule is represented by a dotted line. For small  $M$ , the mean is far from the Bayes prediction rule (green line), and the learned prediction rules have low variance. As  $M$  increases, the mean value comes closer to the Bayes prediction rule, but the variance increases. This is further shown in Figure 3.10 that shows the function  $b_n(x)$  and  $v_n(x)$  for different polynomial orders. For small  $M$ , the bias term is large and the variance term small. The variance increases with the model order, while the bias term decreases. The bias-variance errors are shown in Figure 3.11 as a function of the order of the model.

### 3.4 Regularised Empirical Risk Minimisation

An alternative to considering classes of different capacity is to consider a single (large) class of prediction rules, and penalise more complex prediction rules within this class. Let  $\text{pen} : \mathcal{F} \rightarrow \mathbb{R}_+$  be a **penalty** or **regularisation** function. For a given prediction rule  $h$ ,  $\text{pen}(h)$  is a measure of the complexity of the prediction rule  $h$ ; larger values correspond to more complex prediction rules, while smaller values are less complex prediction rules.

**Definition 21.** Let  $d$  be a dataset. For a loss function  $L$ , a hypothesis class  $\mathcal{H} \subset \mathcal{F}$ , a penalty function  $\text{pen}$  and a **regularisation parameter**  $\lambda \geq 0$ , the regularised empirical risk minimiser is given by

$$\hat{h}^{(d)} = \arg \min_{h \in \mathcal{H}} \left\{ \hat{R}_n(h) + \frac{\lambda}{n} \text{pen}(h) \right\}. \quad (3.8)$$

The empirical risk  $\hat{R}_n(h)$  measures the fit of the prediction rule to the data. More complex prediction rules will typically give a better fit on the data. The second term  $\frac{\lambda}{n} \text{pen}(h)$  penalises the complexity of the prediction rule. The best prediction is the prediction rule achieving the best compromise between the fit to the data and the complexity. The regularisation parameter  $\lambda$  controls this trade-off. If  $\lambda = 0$ , this corresponds to the ERM

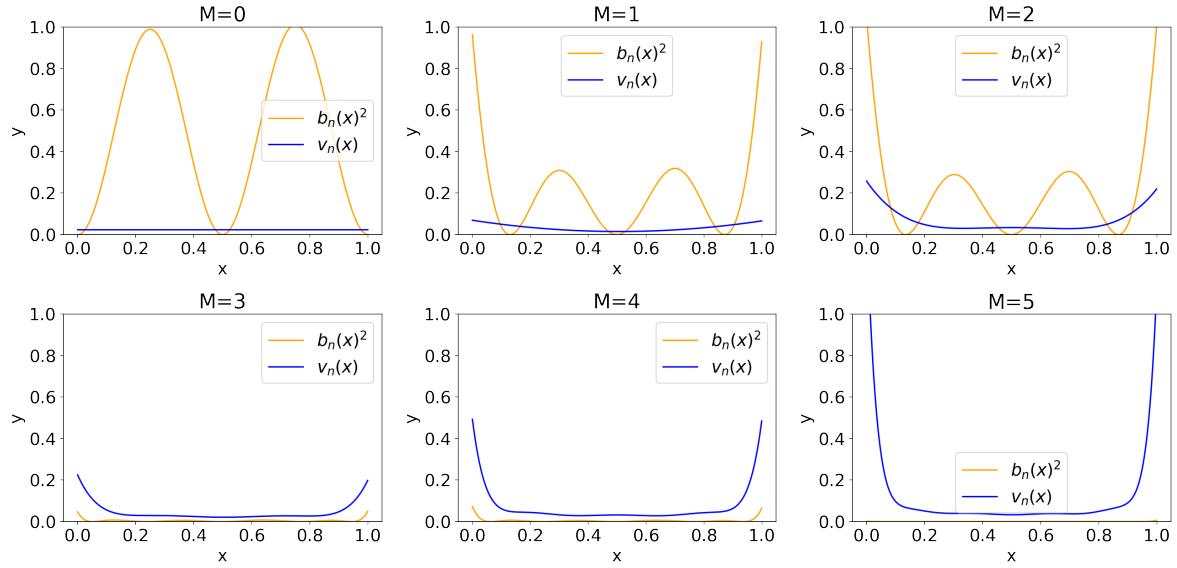


Figure 3.10: Simulated sinusoid example, with  $n = 30$ . Bias and variance as functions of  $x$  for different orders of the polynomial. For small order  $M$ , the bias term dominates, while for large  $M$ , the variance term dominates.

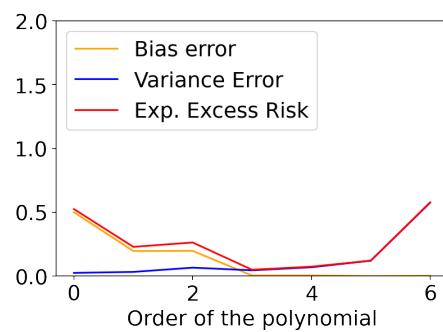


Figure 3.11: Bias-variance error and expected excess risk for the sinus dataset ( $n = 30$ ).

solution: there is no penalty for the complexity, and the learned prediction rule is likely to overfit if the class is large. As  $\lambda \rightarrow \infty$ , a very large penalty is enforced for the model's complexity, and simpler prediction rules will be chosen. In terms of bias-variance decomposition, larger values of  $\lambda$  will lead to higher bias and smaller variance.

If the prediction rules  $h \in \mathcal{H}$  are parameterised by a vector  $\beta = (\beta_0, \dots, \beta_M)^\top \in \mathbb{R}^{M+1}$ , a standard measure of the complexity of the prediction rule is the square of the L2 norm  $\|\beta\|^2 = \sum_{j=0}^M \beta_j^2$  of the vector  $\beta$ . The associated penalty, known as **Tikhonov regularisation**, is therefore

$$\text{pen}(h) = \|\beta\|^2 = \sum_{j=0}^M \beta_j^2.$$

For example, for linear regression under the squared loss, if  $\mathcal{H}$  is the class of linear prediction rules on some transformed inputs  $\phi(x)$ ,

$$\mathcal{H} = \{h(x) = \phi(x)^\top \beta \text{ where } \beta \in \mathbb{R}^{M+1}\}$$

the regularised ERM is  $\hat{h}^{(d)}(x) = \phi(x)^\top \hat{\beta}$  where

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{M+1}} \sum_{i=1}^n (y_i - \phi(x_i)^\top \beta)^2 + \lambda \|\beta\|^2. \quad (3.9)$$

This is the **ridge regression estimator**, which admits a closed-form solution

$$\hat{\beta} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{y}, \quad (3.10)$$

where  $I$  is the identity matrix and  $\Phi$  and  $\mathbf{y}$  are defined in Equation (3.2).

*Proof.*

$$\begin{aligned} \sum_{i=1}^n (y_i - \phi(x_i)^\top \beta)^2 + \lambda \|\beta\|^2 &= \|\mathbf{y} - \Phi \beta\|^2 + \lambda \|\beta\|^2 \\ &= \beta^\top \Phi^\top \Phi \beta - 2(\Phi^\top \mathbf{y})^\top \beta + \lambda \beta^\top \beta + \text{const.} \end{aligned}$$

Differentiating with respect to  $\beta$ , using (3) and (4), we obtain

$$\frac{\partial(\|\mathbf{y} - \Phi \beta\|^2 + \lambda \|\beta\|^2)}{\partial \beta} = 2\Phi^\top \Phi \beta - 2\Phi^\top \mathbf{y} + 2\lambda \beta.$$

Setting the derivative to 0 gives the ridge regression estimator

$$\hat{\beta} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{y}.$$

□

**Example 22.** (sinusoid, continued) For  $n = 10$  observations, we consider the class of linear models with polynomial expansion of order 9, adding a regularisation term  $\lambda$ . The learned prediction rules are shown in Figure 3.12, and the value of the risk and empirical risk are shown in Figure 3.13(middle). As the value of the regularisation parameter  $\lambda$  increases, the learned prediction rule tends to have smaller coefficients, which are shrunk towards 0 (see Figure 3.13(left)). When  $\lambda$  goes to 0, the learned prediction rule tends to the non-regularised solution which overfits the data. For  $n = 30$  observations, we give in Figure 3.13(right).

**Regularisation and sample size.** As the sample size  $n \rightarrow \infty$ , for any fixed  $h$ ,

$$\begin{aligned} \frac{\lambda}{n} \text{pen}(h) &\rightarrow 0 \\ \widehat{R}_n(h) &\rightarrow R(h) \quad \text{almost surely} \end{aligned}$$

The effect of the regularisation therefore vanishes as we have more data.

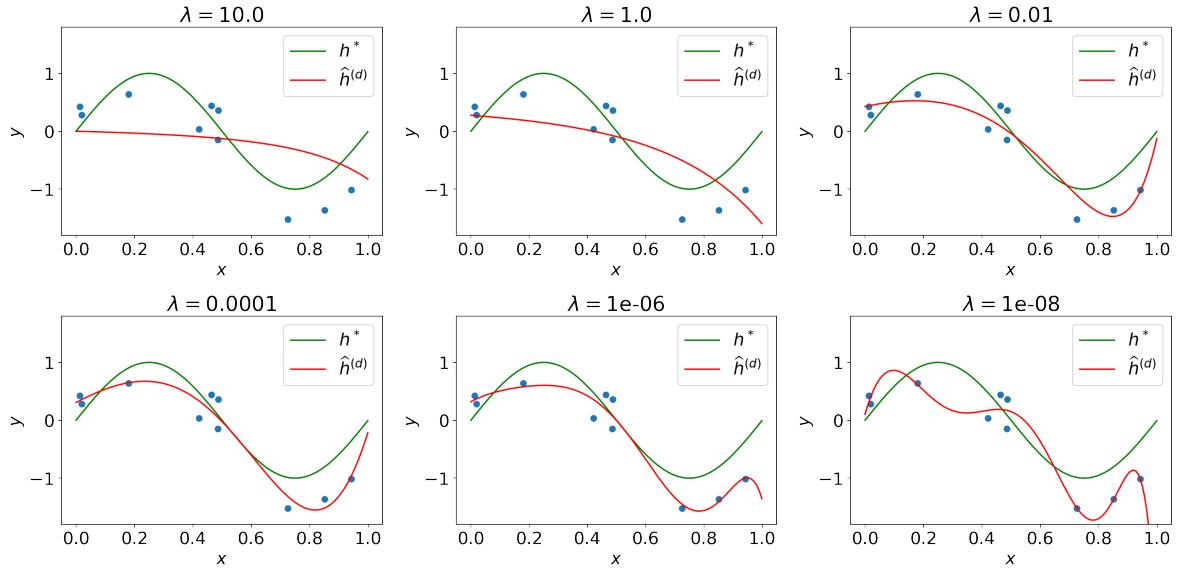


Figure 3.12: Simulated sinusoid example, with  $n = 10$ . Learned prediction rule for different values of the regularisation parameter  $\lambda$ .

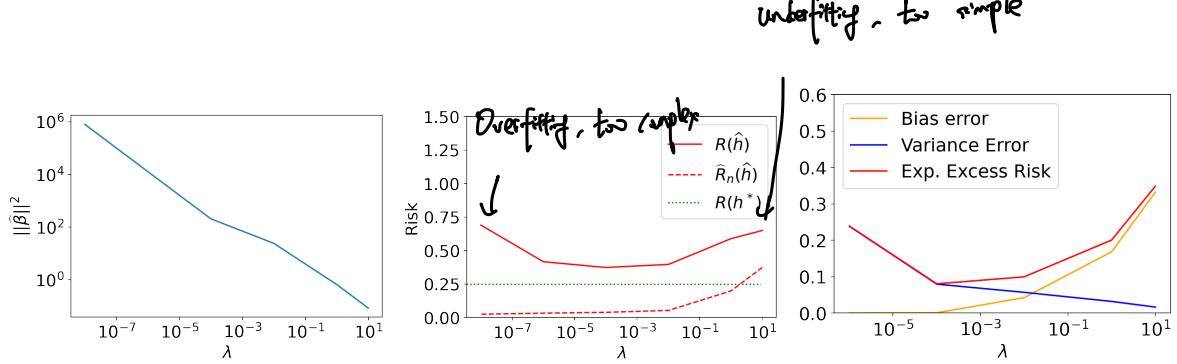


Figure 3.13: (Left) Norm of the estimated parameter  $\beta$  as a function of the regularisation parameter  $\lambda$ , on a log-log plot. (Middle) Risk and empirical risk as a function of the regularisation parameter  $\lambda$ . (Right) Bias vs Variance error as a function of the regularisation parameter  $\lambda$ .

**Other penalisation functions.** Alternative penalisation functions are also commonly used.

- L1 penalty

$$\text{pen}(h) = 2\|\beta\|_1 = 2 \sum_{j=0}^M |\beta_j|$$

- L1+L2 (elastic net)

$$\text{pen}(h) = 2\delta\|\beta\|_1 + (1 - \delta)\|\beta\|_2^2$$

for some  $\delta \in [0, 1]$ .

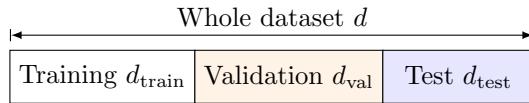
### 3.5 Test, validation and cross-validation

Assume that we obtain a set of  $M_{\max}$  learned prediction rules  $\hat{h}_1^{(d)}, \dots, \hat{h}_{M_{\max}}^{(d)}$ . These prediction rules may for example be obtained using ERM on hypothesis sets  $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots \subset \mathcal{H}_{M_{\max}}$ , of increasing complexity. They may be obtained by using ERM on a large class  $\mathcal{H}$ , with different regularisation parameters  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{M_{\max}}$ . More generally, they may be obtained by using different hyperparameters, learning algorithms, feature expansions, etc.

As we have seen, some learned prediction rules will overfit or underfit. We would like to choose the prediction rule with the smallest generalisation risk, and report this risk, as an estimate of the average error for predicting new observations. There are two caveats:

- The true risk cannot be computed as the true distribution of the data is unknown;
- The empirical risk of the learned prediction rule can be computed but, as we have seen, it largely underestimates the true risk, in particular for complex models.

The idea is to split the dataset in three parts: a training set, a validation set, and a test set.



Denote  $d_{\text{train}}$ ,  $d_{\text{val}}$  and  $d_{\text{test}}$  respectively the training, validation and test set. As before, we assume that the training, validation and test set are realisations from  $D_{\text{train}}$ ,  $D_{\text{val}}$  and  $D_{\text{test}}$ . Denote  $\hat{R}^{(d_{\text{train}})}(h)$ ,  $\hat{R}^{(d_{\text{val}})}(h)$  and  $\hat{R}^{(d_{\text{test}})}(h)$  the empirical risk of a prediction rule  $h$  on the training, validation and training sets. The procedure is then as follows.

#### Algorithm 1 Training-Validation-Test procedure

1. Training: For each  $j = 1, \dots, M_{\max}$ , estimate the learned prediction rule  $\hat{h}_j^{(d_{\text{train}})}$  using the training set  $d_{\text{train}}$ . For instance, when considering the ERM on a class  $\mathcal{H}_j$

$$\hat{h}_j^{(d_{\text{train}})} = \arg \min_{h \in \mathcal{H}_j} \hat{R}^{(d_{\text{train}})}(h).$$

2. Validation (model selection):

- For each  $j$ , calculate  $\hat{R}^{(d_{\text{val}})}(\hat{h}_j^{(d_{\text{train}})})$  and report the best learner  $\widehat{M} = \arg \min_j \hat{R}^{(d_{\text{val}})}(\hat{h}_j^{(d_{\text{train}})})$
- Re-train the prediction rule on both the training and validation set for the chosen class  $\widehat{M}$ . For instance, for ERM on a class  $\mathcal{H}_{\widehat{M}}$

$$\hat{h}^{(d_{\text{train}}, d_{\text{val}})} = \arg \min_{h \in \mathcal{H}_{\widehat{M}}} R^{(d_{\text{train}}, d_{\text{val}})}(h)$$

3. Test: Approximate the generalisation error with  $\hat{R}^{(d_{\text{test}})}(\hat{h}^{(d_{\text{train}}, d_{\text{val}})})$ .

Some general comments:

- In Step 1, as we use the training set  $d_{\text{train}}$  to estimate the  $\hat{h}_j^{(d_{\text{train}})}$ , we cannot use  $\hat{R}^{(d_{\text{train}})}(\hat{h}_j^{(d_{\text{train}})})$  as an estimate for  $R(\hat{h}_j^{(d_{\text{train}})})$
- We therefore use in step 2a) another dataset (the validation set), that has not been used for training. Note that, for any  $j$ , and any training set  $d_{\text{train}}$ ,

$$\mathbb{E}[\hat{R}^{(d_{\text{val}})}(\hat{h}_j^{(d_{\text{train}})}) | D_{\text{train}} = d_{\text{train}}] = R(\hat{h}_j^{(d_{\text{train}})})$$

where the expectation is taken over the validation set. It therefore provides an unbiased estimate of the true risk of the  $j$ th learned prediction rule.

- After selecting the best class  $\widehat{M}$ , we could use  $\widehat{h}_{\widehat{M}}^{(d_{\text{train}})}$  as the estimated prediction rule. However, to reduce the estimation error, it is better to re-train on both the training and validation sets for this model.
- The training and validation datasets having been used to estimate  $\widehat{h}^{(d_{\text{train}}, d_{\text{val}})}$ , one has to calculate the empirical on a third dataset (the test set). We have

$$\mathbb{E}[\widehat{R}^{(D_{\text{test}})}(\widehat{h}^{(D_{\text{train}}, D_{\text{val}})}) \mid D_{\text{train}} = d_{\text{train}}, D_{\text{val}} = d_{\text{val}}] = R(\widehat{h}^{(d_{\text{train}}, d_{\text{val}})})$$

where the expectation is taken over the test set. It thus provides an unbiased estimate of the true risk of the learned prediction rule.

How to split the dataset? Any data in the validation and test set is not used to initially train the different models. If the training set is too small, this may lead to large estimation errors for the learned prediction rules. Alternatively, if the validation set or test sets are small, this leads to a large error on estimated risks. In the case where the dataset is small or the learners need a lot of data, an alternative is to use cross-validation.

**Cross-validation.** The idea of  $T$ -fold cross-validation is to split the dataset into a training/validation set and a test set, then to further split the training/validation set into  $T$  folds.

1. Training: For each fold  $t = 1, \dots, T$

- Use fold  $t$  as a validation set, and the rest as a training set to train each learner  $j$ ; denote  $\widehat{h}_j^{(d_{\text{train},t})}$  the  $j$ th learned prediction rule, and

$$\widehat{R}^{(d_{\text{val},t})}(\widehat{h}_j^{(d_{\text{train},t})})$$

its estimated risk

2. Validation:

- a) Choose the learner that minimises the estimated risk average over the folds

$$\widehat{M} = \arg \min_j \frac{1}{T} \sum_{t=1}^T \widehat{R}^{(d_{\text{val},t})}(\widehat{h}_j^{(d_{\text{train},t})})$$

- b) Re-train the prediction rule on the whole training/validation set for the chosen learner  $\widehat{M}$ .

3. Test: Approximate the generalisation error with  $\widehat{R}^{(d_{\text{test}})}(\widehat{h}^{(d_{\text{train}/\text{val}})})$ .

$T=4$

Whole dataset				
d <sub>val,1</sub>		d <sub>train,1</sub>		
Val	Train	Train	Train	Test
Train	Val	Train	Train	Test
Train	Train	Val	Train	Test
Train	Train	Train	Val	Test

A special case of cross-validation is **leave-one-out** cross-validation, where we use one data item per fold, that is  $T$  equals the number of observations in training/validation set. Cross-validation can be computationally intensive, as the computational cost is multiplied by the number of folds  $T$ .

### 3.6 Binary Classification: Performance evaluation and ROC curve

We focus here on the binary classification case. Let  $y$  be the true class, and  $h(x)$  be the predicted class. The following table summarises the different configurations.

~~x~~

		$h(x) = -1$	$h(x) = 1$
$y = -1$	true negative	false positive	
$y = 1$	false negative	true positive	

$$TN = \sum_{i=1}^N \mathbb{1}_{y_i = h(x_i) = -1}$$

The matrix counting the number of true/false positive/negative over a dataset  $d$  is called a **confusion matrix**. We denote  $TN, FN, FP, TP$  respectively the number of true negative, false negative, false positive and true positive over the dataset  $d$  for the classifier  $h$ . So far, as a measure of the quality of a classifier  $h$ , we have focused on the (population) risk under a 0-1 loss, which is the misclassification error, or error rate

$$R(h) = \Pr(Y \neq h(X)).$$

The empirical misclassification error is

$$\hat{R}_n^{(d)}(h) = \frac{FN + FP}{FN + FP + TP + TN}.$$

For binary classification, other quantities are often considered to assess the performance of a classifier. We list below the population and empirical versions.

Name	Population	Empirical
Misclassification error/Error rate	$\Pr(Y \neq h(X))$	$(FP + FN)/(TP + TN + FP + FN)$
Accuracy (1-Error rate)	$\Pr(Y = h(X))$	$(TP + TN)/(TP + TN + FP + FN)$
Sensitivity/Recall/True positive rate	$\Pr(h(X) = 1   Y = 1)$	$TP/(TP + FN)$
Specificity/True negative rate	$\Pr(h(X) = -1   Y = -1)$	$TN/(TN + FP)$
False positive rate (1-specificity)	$\Pr(h(X) = 1   Y = -1)$	$FP/(TN + FP)$
Precision	$\Pr(Y = 1   h(X) = 1)$	$TP/(TP + FP)$

**Weighted loss.** Two possible types of error may occur:

- False positive ( $y = -1, h(x) = 1$ ), also called Type I error,
- False negative ( $y = 1, h(x) = -1$ ), also called Type II error.

The 0 – 1 loss assigns an equal cost of 1 to each type of error

$$\begin{array}{c|cc} & h(x) = -1 & h(x) = 1 \\ \hline y = -1 & L(y, h(x)) = 0 & L(y, h(x)) = 1 \\ y = 1 & L(y, h(x)) = 1 & L(y, h(x)) = 0 \end{array}$$

In some cases, it is more sensible to assign different costs to the different types of error. For example, when predicting if a patient has a given disease, we may assign a higher cost to a false negative error than a false positive. We consider in this case a loss

$$L(y, h(x)) = a\mathbb{1}_{y=-1, h(x)=1} + b\mathbb{1}_{y=1, h(x)=-1}$$

where  $a, b > 0$  are respectively the type I and type II cost. The Bayes classifier under this cost is (see problem sheet)

$$h_t^*(x) = \begin{cases} 1 & \text{if } \eta(x) \geq t := \frac{a}{a+b} \\ -1 & \text{otherwise} \end{cases}$$

where  $\eta(x) = \Pr(Y = 1 | X = x)$ . Note that the Bayes classifier only depends on  $a$  and  $b$  through the ratio  $t = a/(a + b)$ , the value  $t = 1/2$  corresponding to the Bayes classifier under the 0 – 1 loss. A large value  $t$  will allow more false negative and less false positive; a small value allows for more false positive and less false negative.

**ROC curve and AUC.** For a given weighted loss with ratio  $t$ , let  $\hat{h}_t$  be some plug-in (generative or discriminative) classifier, defined by

$$\hat{h}_t(x) = \begin{cases} 1 & \text{if } \hat{\eta}(x) \geq t \\ -1 & \text{otherwise} \end{cases}$$

where  $\hat{\eta}(x)$  is the estimate of  $\eta(x) = \Pr(Y = 1 | X = x)$ , and does not depend on the choice of the threshold  $t$ .

How can we assess the performance of the family of plug-in classifiers  $(\hat{h}_t)$  for  $t \in [0, 1]$ ? A classical strategy is the **ROC curve** and the **Area Under the Curve (AUC)**. For a threshold  $t \in [0, 1]$ , denote  $\beta(t) = \Pr(\hat{h}_t(X) = 1 | Y = 1)$  the true positive rate (TPR) of  $\hat{h}_t$  and  $\alpha(t) = \Pr(\hat{h}_t(X) = 1 | Y = -1)$  its false positive rate (FPR).  $\alpha$  and  $\beta$  are monotone functions of  $t$ , with  $\alpha(t) = \beta(t) = 0$  if  $t = 1$ , and  $\alpha(t) = \beta(t) = 1$  if  $t = 0$ . A good classifier will have a large true positive rate for a small false positive rate. To quantify this over a range of values, the ROC curve plots the TPR  $\beta(t)$  vs the FPR  $\alpha(t)$  for  $t$  ranging from 1 to 0. The area under the ROC curve, called AUC and defined as

$$AUC = \int_1^0 \beta(t)d\alpha(t),$$

is used as a measure of the performance of the family of plug-in classifiers over the whole range of loss functions. This is illustrated in Figure 3.14. In practice the population TPR and FPR cannot be computed analytically, and one reports their empirical approximation over a test set.

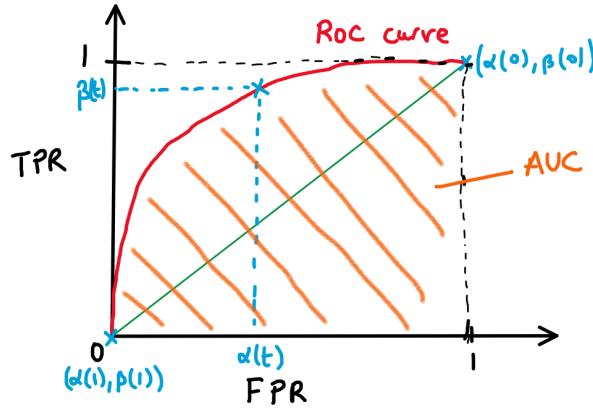


Figure 3.14: Illustration of the ROC curve and Area Under the Curve. Each point of the ROC curve corresponds to the value of the False positive rate vs the true positive rate of a plug-in classifier for a threshold  $t$ .

**Probabilistic Interpretation of the AUC.** Consider a plug-in classifier with estimate  $\hat{\eta}(x)$ , obtained from a dataset  $d$ , which is a realisation from a random sample  $D$ . For two random variables  $(\tilde{X}_0, \tilde{Y}_0)$  and  $(\tilde{X}_1, \tilde{Y}_1)$ , mutually independent and independent from  $D$ , with the same distribution as  $(X, Y)$ , we have

$$\mathcal{D} \quad AUC = \Pr(\hat{\eta}(\tilde{X}_1) \geq \hat{\eta}(\tilde{X}_0) \mid \tilde{Y}_1 = 1, \tilde{Y}_0 = -1, D = d).$$

That is the AUC is the probability that, for two independent samples from class 1 and  $-1$ , the discriminant score of the sample from class 1 is greater than the discriminant score of the sample from class  $-1$ .

*Proof.* The conditioning on  $D = d$  is implicit in what follows. Denote  $F_{-1}(s)$  the cdf of the random variable  $\hat{\eta}(X) \in [0, 1]$  given  $Y = -1$ , and  $F_1(s)$  the cdf of  $\hat{\eta}(X)$  given  $Y = 1$ :

$$\begin{aligned} F_1(s) &= \Pr(\hat{\eta}(X) \leq s \mid Y = 1) \\ F_{-1}(s) &= \Pr(\hat{\eta}(X) \leq s \mid Y = -1). \end{aligned}$$

Let  $f_1 = F'_1$  and  $f_{-1} = F'_{-1}$  be their pdf. For a threshold  $t = a/(a+b)$ , both the (population) true positive rate  $\beta(t)$  and false positive rate  $\alpha(t)$  of the plug-in classifier  $\hat{h}_t$  can be expressed as

$$\beta(t) = \Pr(\hat{\eta}(X) \geq t \mid Y = 1) = 1 - F_1(t) \tag{3.11}$$

$$\alpha(t) = \Pr(\hat{\eta}(X) \geq t \mid Y = -1) = 1 - F_{-1}(t) \tag{3.12}$$

The Area Under the ROC Curve, denoted AUC is

$$\begin{aligned} \int_1^0 \beta(t) d\alpha(t) &= \int_0^1 (1 - F_1(t)) f_{-1}(t) dt = \mathbb{E}[(1 - F_1(\hat{\eta}(\tilde{X}_0))) \mid \tilde{Y}_0 = -1] \\ &= \mathbb{E}[\Pr(\hat{\eta}(\tilde{X}_1) \geq \hat{\eta}(\tilde{X}_0) \mid \tilde{Y}_1 = 1) \mid \tilde{Y}_0 = -1] \\ &= \Pr(\hat{\eta}(\tilde{X}_1) \geq \hat{\eta}(\tilde{X}_0) \mid \tilde{Y}_1 = 1, \tilde{Y}_0 = -1) \end{aligned}$$

where  $(\tilde{X}_0, \tilde{Y}_0)$  and  $(\tilde{X}_1, \tilde{Y}_1)$  are iid with the same distribution as  $(X, Y)$ .  $\square$

# 4 — Optimisation for machine learning

Let  $\theta \in \mathbb{R}^p$ . Consider either

- ERM with a class of prediction rules  $h_\theta$  parameterised by a parameter  $\theta$  or
- a conditional plug-in method, where the conditional distribution  $f_\theta(y|x)$  is parameterised by  $\theta$  or
- a generative method, where the joint distribution  $\pi_\theta(x, y)$  is parameterised by  $\theta$ .

For the plug-in approach, we assume that the parameter  $\theta$  is fitted using maximum likelihood, potentially with the addition of a regularisation term.

In all three cases, we aim at minimising an objective function  $J : \mathbb{R}^p \rightarrow \mathbb{R}_+$  of the form

$$\underline{J(\theta) = \sum_{i=1}^n J_i(\theta) + J_0(\theta)} \quad (4.1)$$

where  $J_0$  is a penalisation/regularisation term and each term  $J_i$  is associated to an example  $(x_i, y_i)$ :

- \*
- For empirical risk minimisation,  $J_i(\theta) = \frac{1}{n} L(y_i, h_\theta(x_i))$
  - For a conditional plug-in approach,  $J_i(\theta) = \log(f_\theta(y_i|x_i))$
  - For a generative plug-in approach,  $J_i(\theta) = \log(\pi_\theta(x_i, y_i))$ .

To keep the notations simple, in this chapter, we consider an ERM objective function, but the algorithms apply equally to estimate the parameters of plug-in methods. For ERM, the objective function  $J$  is therefore of the form

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(y_i, h_\theta(x_i)) + \underbrace{\frac{\lambda}{n} \text{pen}(h_\theta)}_{\text{penalty}}. \quad (4.2)$$

Assuming it exists, the global minimum to (4.2) may not be available in closed-form and one has to resort to some numerical method. Even if it is available in closed-form, computing it may be impractical for a large number  $n$  of observations and/or a large number of parameters  $p$ .

For example, for empirical risk minimisation with the class of linear functions with input expansion  $\phi : \mathcal{X} \rightarrow \mathbb{R}^p$ , and L2 regularisation with regularisation parameter  $\lambda$ , the objective function is

$$J(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \phi(x_i)^\top \beta)^2 + \frac{\lambda}{n} \|\beta\|^2. \quad (4.3)$$

As shown in Section 3.4, the estimated parameters of the learned prediction rule are computed as

$$\hat{\beta} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{y} \quad \boxed{\Phi = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \quad \text{NP}}$$

where  $\Phi$  is an  $n \times p$  matrix,  $\mathbf{y}$  is a vector of length  $n$ , and  $\Phi^\top \Phi + \lambda I$  is a  $p \times p$  matrix. The number of computations is of order

- $O(np^2)$  for computing  $\Phi^\top \Phi + \lambda I$  and  $\Phi^\top \mathbf{y}$ .
- $O(p^3)$  (using e.g. Gauss-Jordan elimination) for solving the linear system

$$(\Phi^\top \Phi + \lambda I)\beta = \Phi^\top \mathbf{y}. \quad (4.4)$$

This becomes impractical for large  $p$  and/or  $n$ . For other methods, such as logistic regression or neural networks, there is no closed-formed solution. Standard numerical approaches are gradient descent (GD) and stochastic gradient descent (SGD), which are reviewed in this chapter.

**Notations.** For a vector  $\theta = (\theta_1, \dots, \theta_p)^\top$ , a twice differentiable function  $J : \mathbb{R}^p \rightarrow \mathbb{R}_+$ , let  $\nabla_\theta J(\theta)$  be the gradient of  $J$  with respect to  $\theta$ , and  $\nabla_\theta^2 J(\theta)$  be the Hessian of  $J$  with respect to  $\theta$ , defined as

$$\nabla_\theta J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_p} \end{pmatrix}, \quad \nabla_\theta^2 J(\theta) = \begin{pmatrix} \frac{\partial^2 J(\theta)}{\partial \theta_1 \partial \theta_1} & \cdots & \frac{\partial^2 J(\theta)}{\partial \theta_1 \partial \theta_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 J(\theta)}{\partial \theta_p \partial \theta_1} & \cdots & \frac{\partial^2 J(\theta)}{\partial \theta_p \partial \theta_p} \end{pmatrix}.$$

A symmetric real-valued  $p$ -by- $p$  matrix  $H$  is said to be positive semi-definite, iff

$$\underline{z^\top H z \geq 0}$$

for all  $z \in \mathbb{R}^p$ . We use the notation  $H \succeq 0$  if  $H$  is positive semi-definite.

## 4.1 Convex functions

Let  $J : \mathbb{R}^p \rightarrow \mathbb{R}$ . We review a few definitions and basic properties of convex functions.

**Definition 23 (Convex function).** A function  $J$  is said to be convex if, for all  $u, v \in \mathbb{R}^p$ ,  $\alpha \in [0, 1]$ ,

$$\underline{J(\alpha u + (1 - \alpha)v) \leq \alpha J(u) + (1 - \alpha)J(v)}$$

Examples of convex and non-convex functions are given in Figure 4.1. Here are some examples of convex functions:

- Univariate:  $\theta^2$ ,  $\exp(-\theta)$ ,  $\log(1 + \exp(-\theta))$ ,  $\max(0, 1 - \theta)$
- Affine functions:  $A\theta + b$
- Quadratic functions:  $\theta^\top H\theta$  where  $H \succeq 0$

A key property of convex functions is the following, which induces that the convergence of an algorithm to a local minimum implies the convergence to a global minimum.

**Proposition 24.** If the function  $J$  is convex, all local minima are also global minima.

For differentiable and twice differentiable convex functions, we have the following characterisations and properties.

**Definition 25 (Differentiable convex function).** A differentiable function  $J$  is convex iff, for any  $u, v$

$$\underline{J(u) \geq J(v) + \nabla J(v)^\top(u - v)}$$

That is, the function is above the tangent at  $v$ .

**Proposition 26.** Let  $J$  be a differentiable convex function. Any stationary point  $u$  of  $J$ , that is such that

$$\underline{\nabla_u J(u) = 0}$$

is also a global minimum.

**Definition 27 (Twice differentiable convex function).** A twice differentiable function  $J$  is convex iff

$$\underline{\nabla_u^2 J(u) \succeq 0}$$

for all  $u \in \mathbb{R}^p$ .

Convex functions can be combined in a number of ways:

- **Nonnegative linear combination of convex functions is convex.** Let  $J_1$  and  $J_2$  be convex functions. Then

$$J = \alpha_1 J_1 + \alpha_2 J_2$$

is also convex for any  $\alpha_1, \alpha_2 \geq 0$ .

- **Affine composition of convex functions is convex.** If  $g$  is convex, then  $J(\theta) = g(A\theta + b)$  is convex. It follows from the first property above that, if for any  $(x_i, y_i)$ , the function  $J_i : \theta \rightarrow L(y_i, h_\theta(x_i))$  is a convex function (of  $\theta$ ) and the regularisation function pen is convex, then the objective function (4.2) is also a convex function, and all local minima are global minima. This is the case for ridge regression, as shown below. The objective functions of logistic regression and of the other linear classifiers discussed in Chapter 5 are also all convex. Other popular machine learning methods, such as support vector machines (not covered in this course) also have a convex objective function. Note that the composition of convex functions is in general non-convex, and the objective function of (deep) neural networks, covered later in the course, is non-convex.

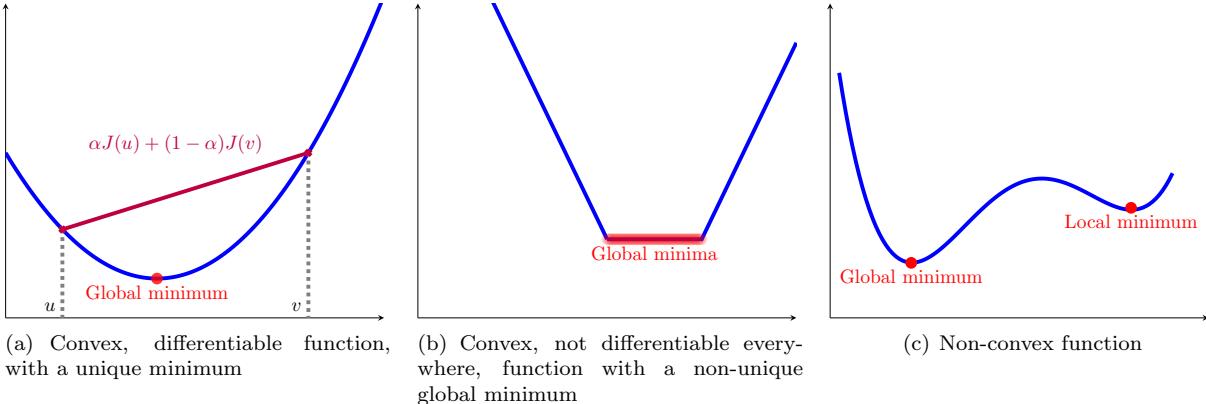


Figure 4.1: Examples of convex and non-convex functions.



**Ridge regression.** The objective function (4.3) is convex. The gradient and Hessian are

$$\begin{aligned} J(\beta) &= \frac{1}{n} \sum_{i=1}^n (y_i - \Phi^\top \beta)^2 + \frac{\lambda}{n} \|\beta\|^2 \\ \nabla_\beta J(\beta) &= \frac{1}{n} (2\Phi^\top \Phi \beta - 2\Phi^\top y + 2\lambda\beta) \\ \nabla_\beta^2 J(\beta) &= \frac{2}{n} (\Phi^\top \Phi + \lambda I). \end{aligned}$$

In this case, the Hessian does not depend on the parameter  $\beta$ , and is positive semi-definite for any  $\lambda \geq 0$ .

## 4.2 Gradient Descent

Consider the objective function (4.2). Denote

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(y_i, h_\theta(x_i)) + \frac{\lambda}{n} \text{pen}(h_\theta)$$

$$\nabla_\theta J(\theta) = \frac{1}{n} \left( \sum_{i=1}^n \nabla_\theta L(y_i, h_\theta(x_i)) + \lambda \nabla_\theta \text{pen}(h_\theta) \right) \quad (4.5)$$

the gradient of  $J$  with respect to  $\theta$ . Let  $\eta > 0$ . The gradient descent algorithm proceeds as follows.

### Algorithm 2 Gradient Descent

$\eta$ : learning rate / step size.

- Initialise  $\theta^{(0)}$  and set  $t = 0$ ;
- Repeat until convergence
  1. Compute the gradient  $\nabla_\theta J(\theta^{(t)})$
  2. Update the parameters

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_\theta J(\theta^{(t)})$$

3. Set  $t \leftarrow t + 1$

An illustration of the iterations of the gradient descent algorithm is given in Figure 4.2(a). The negative gradient is going in the direction that decreases the function  $J$ . Large values of the gradient will induce large changes in the value of the parameter. As the parameter becomes closer to a stationary point of  $J$ , the gradient takes smaller values, and the parameter changes more slowly. For convex functions (under some additional assumptions omitted here), the algorithm will stop at a point close to the minimum solution, independent of the starting point. For non-convex functions, the algorithm will attain different local minima depending on the initialisation.

**Learning rate.** The parameter  $\eta$  is known as the step size, or learning rate. Its value is critical to the good behaviour of the gradient descent algorithm. Too large values may lead the algorithm to have a zigzagging behaviour, or even to diverge, whereas too small values will lead to slow convergence. See the figure 4.2 for an illustration. This parameter can also be updated adaptively.

**Interpretation of gradient descent via a quadratic approximation, and Newton-Raphson algorithm.** Assume that  $J$  is twice differentiable. Using a second-order Taylor expansion of the objective function

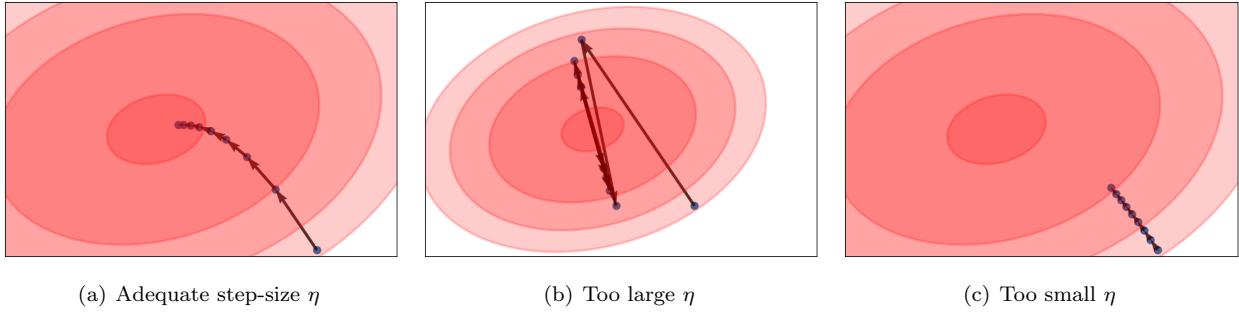


Figure 4.2: Illustration of the first iterations of the gradient descent algorithm for different step-sizes  $\eta$ . Ellipses represent contour of constant value of the objective function.

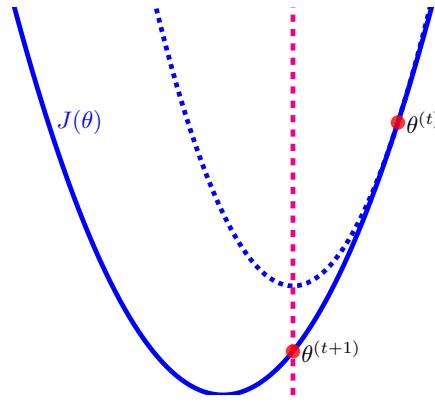


Figure 4.3: Gradient descent as minimisation of a quadratic function.

$J$  around a point  $\theta \in \mathbb{R}^p$ , we have, for  $\delta \in \mathbb{R}^p$ ,

$$J(\theta + \delta) \simeq J(\theta) + \nabla_\theta J(\theta)^\top \delta + \frac{1}{2} \delta^\top \nabla_\theta^2 J(\theta) \delta. \quad (4.6)$$

Replacing the Hessian matrix  $\nabla_\theta^2 J(\theta)$  with the diagonal  $p$ -by- $p$  matrix  $\frac{1}{\eta} I$  gives the approximation

$$J(\theta + \delta) \simeq J(\theta) + \nabla_\theta J(\theta)^\top \delta + \frac{1}{2\eta} \|\delta\|^2. \quad (4.7)$$

At a given point  $\theta$ , we use the above quadratic function as an approximation to  $J(\theta + \delta)$ . Taking the derivative with respect to  $\delta$ , we obtain

$$\nabla_\theta J(\theta) + \frac{1}{\eta} \delta = 0,$$

hence the minimum is achieved for  $\delta = -\eta \nabla_\theta J(\theta)$ , which corresponds to the gradient descent update. See Figure 4.3 for an illustration. Instead of taking approximating the Hessian, we can use directly the second-order approximation (4.6). Differentiating with respect to  $\delta$  and setting to 0, we obtain

$$\nabla_\theta J(\theta) + \nabla_\theta^2 J(\theta) \delta = 0$$

hence a minimum at

$$\delta = -(\nabla_\theta^2 J(\theta))^{-1} \nabla_\theta J(\theta).$$

The corresponding update is therefore  $\theta^{(t+1)} = \theta^{(t)} - (\nabla_\theta^2 J(\theta^{(t)})^{-1} \nabla_\theta J(\theta^{(t)})$ , and the iterative algorithm using this update is known as the **Newton-Raphson** gradient descent algorithm. Note that the Hessian is a  $p$ -by- $p$  matrix, which may be computationally too expensive to compute for large  $p$ .

**Gradient descent for linear regression.** For instance, for ordinary linear regression (no regularisation), we have

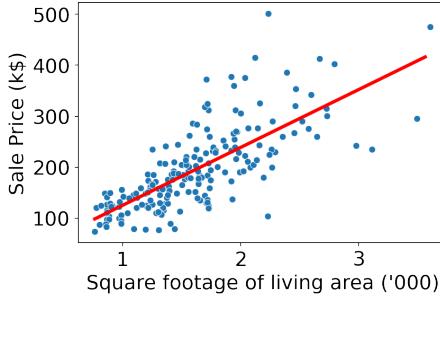
$$nJ(\beta) = (\mathbf{y} - \Phi\beta)^\top (\mathbf{y} - \Phi\beta). \quad (4.8)$$

**Algorithm 3** Newton-Raphson algorithm

- Initialise  $\theta^{(0)}$  and set  $t = 0$ ;
- Repeat until convergence
  1. Update the parameters

$$\theta^{(t+1)} = \theta^{(t)} - (\nabla_{\theta}^2 J(\theta^{(t)}))^{-1} \nabla_{\theta} J(\theta^{(t)})$$

2. Set  $t \leftarrow t + 1$



$$\begin{aligned}
 n J(\beta) &= (\mathbf{y} - \Phi \beta)^T (\mathbf{y} - \Phi \beta) \\
 n \frac{\partial J(\beta)}{\partial \beta_1} &= 2(\mathbf{y}_1 - \Phi \beta_1) \beta_1 - \Phi^T \mathbf{y}_1 \\
 n \nabla_{\beta_1} J(\beta) &= 2 \Phi^T (\Phi \beta - \mathbf{y})
 \end{aligned}$$

Figure 4.4: House pricing data and learned prediction rule.

The gradient is expressed as

$$n \nabla_{\beta} J(\beta) = 2 \Phi^T (\Phi \beta - \mathbf{y})$$

leading to the parameter update

$$\beta^{(t+1)} = \beta^{(t)} - \frac{2\eta}{n} \Phi^T (\Phi \beta^{(t)} - \mathbf{y}).$$

The initial cost of gradient descent is  $O(np^2)$  to compute  $\Phi^T \Phi$ . Then each iteration has a computational cost of  $O(np)$ . If  $n$  and  $T$ , the total number of gradient descent steps, are both small compared to the number of parameters  $p$ , the overall computational cost of the gradient descent algorithm is lower than that of solving the .

**Input/Feature normalisation.** We can rewrite the gradient step a bit differently to gain some insights on the gradient descent update, and on its limitations:

$$\begin{aligned}
 \beta^{(t+1)} &= \beta^{(t)} - \frac{2\eta}{n} (\Phi^T \Phi \beta^{(t)} - \Phi^T \mathbf{y}) \\
 &= \beta^{(t)} + \frac{2\eta}{n} \Phi^T \Phi (\hat{\beta} - \beta^{(t)})
 \end{aligned}$$

where

$$\hat{\beta} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

is the empirical risk minimiser. If  $\Phi^T \Phi = c \times I$  for some constant  $c > 0$ , then a gradient descent update will go in the direction of  $\hat{\beta} - \beta^{(t)}$ , and will converge quickly to  $\hat{\beta}$ . Otherwise, the update will have a zigzagging behaviour. A way to have a better conditioned problem is to apply some transformation to the data. Standardising the different dimensions of the transformed inputs  $\phi(x_i)$  so that they have zero mean and unit variance for example leads to a more spherical  $\Phi^T \Phi$  and a faster convergence of the GD algorithm. This is illustrated in the example below.

**Example 28.** We consider again the house sale prices dataset introduced in Section 1.4.3 ( $n = 200, p = 1$ ). We first consider a linear prediction rule with transformation  $\phi(x_i) = (1 \ x_i)^T$ . Let  $\hat{\beta} \in \mathbb{R}^2$  be the ordinary least square estimate. The learned prediction rule  $\hat{h}(x) = x^T \hat{\beta}$  is shown in Figure 4.4, together with the data. The estimate  $\hat{\beta}$  can be found in closed-form, but we are going to estimate it using gradient descent, starting from  $\beta^{(0)} = (0, 0)^T$ . The first iterates, and the result after 200 iterations are represented in Figure 4.5. As can be seen, the contours of constant value of the empirical risk have an ellipsoid shape and the negative gradients do not point towards the minimum of the objective function. The algorithm takes a large number of iterations to converge, and is zigzagging. Increasing or decreasing the learning rate  $\eta$  would not resolve this issue. We can obtain a much faster convergence by standardising the inputs. Consider now the transformed inputs  $\phi(x_i) = (1, \frac{x_i - \bar{x}}{\sigma})^T$  where  $\bar{x}$  and  $\sigma$  respectively denote the sample mean and standard deviation of  $(x_1, \dots, x_n)$ . Starting from the

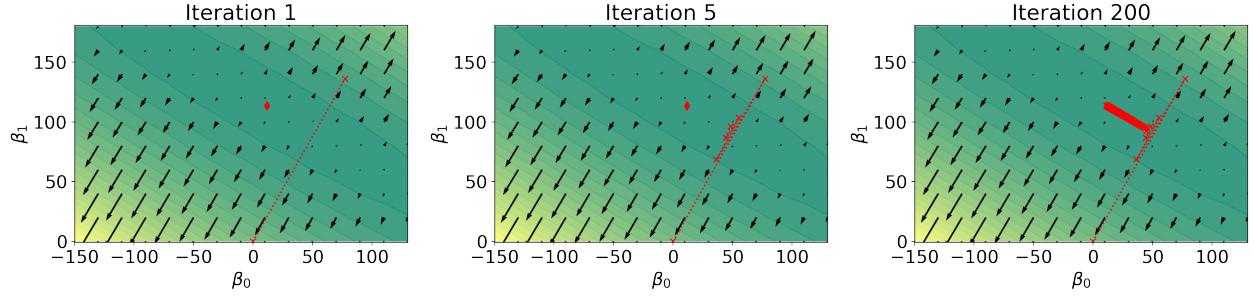


Figure 4.5: Contours of the objective function (unnormalised case), minimum of the objective function (red diamond), and gradient descent updates (crosses). The gradients are represented by arrows.

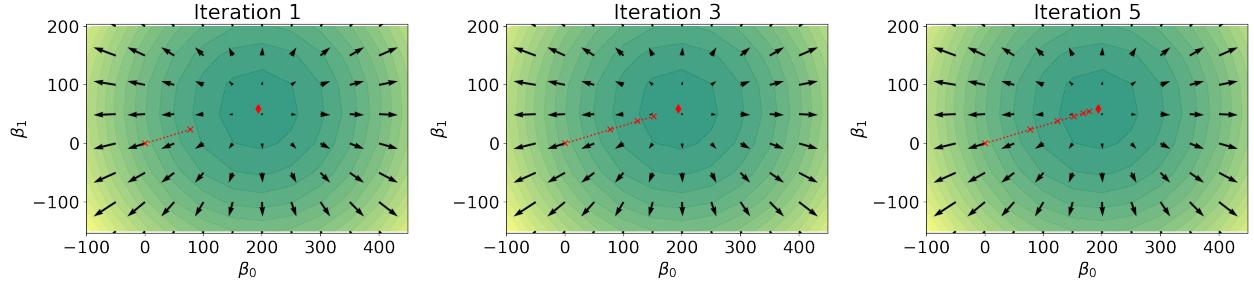


Figure 4.6: Contours of the objective function (normalised case), minimum of the objective function (red diamond), and gradient descent updates (crosses). The gradients are represented by black arrows.

same initialisation, the first iterates are given in Figure 4.6. The contours of the objective function, in this particular case, are now spherical, and the negative gradients point towards the minimum. Using a not too large learning rate, gradient descent converges in a few iterations.

### 4.3 Stochastic Gradient Descent

Stochastic gradient descent (and its variants) is certainly the most popular algorithm for large-scale (large  $n$  and/or  $p$ ) machine learning. It is in particular the default algorithm to learn deep learning models.

Gradient descent requires to evaluate the gradient (4.5) at each iteration, which typically scales linearly in the number  $n$  of observations. This is prohibitive for very large datasets. Stochastic gradient descent replaces the gradient by an unbiased estimator of the gradient, obtained by using a subset of the data (called a **mini-batch**) at each iteration. Let  $(\eta_t)_{t=0,1,\dots}$  be a monotone decreasing sequence, and  $1 \leq n_b \leq n$  be the **mini-batch size**. The algorithm proceeds as follows.




---

#### Algorithm 4 Stochastic Gradient Descent

---

- Initialise  $\theta^{(0)}$  and set  $t = 0$ ;
- Repeat until convergence
  1. Randomly select  $n_b$  observations from the dataset  $d$ ; denote them  $(\tilde{x}_i^{(t)}, \tilde{y}_i^{(t)})_{i=1,\dots,n_b}$
  2. Compute the gradient estimate  $\nabla_\theta \tilde{J}^{(t)}(\theta^{(t)})$  where

$$\nabla_\theta \tilde{J}^{(t)}(\theta) = \frac{1}{n_b} \sum_{i=1}^{n_b} \nabla_\theta L(\tilde{y}_i^{(t)}, h_\theta(\tilde{x}_i^{(t)})) + \frac{\lambda}{n} \nabla_\theta \text{pen}(h_\theta)$$

3. Update the parameters

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \tilde{J}^{(t)}(\theta^{(t)})$$

4. Set  $t \leftarrow t + 1$
- 

For any  $\theta$ , the mini-batch gradient  $\nabla_\theta \tilde{J}^{(t)}(\theta)$  is an unbiased estimator of the true gradient  $\nabla_\theta J(\theta)$ .

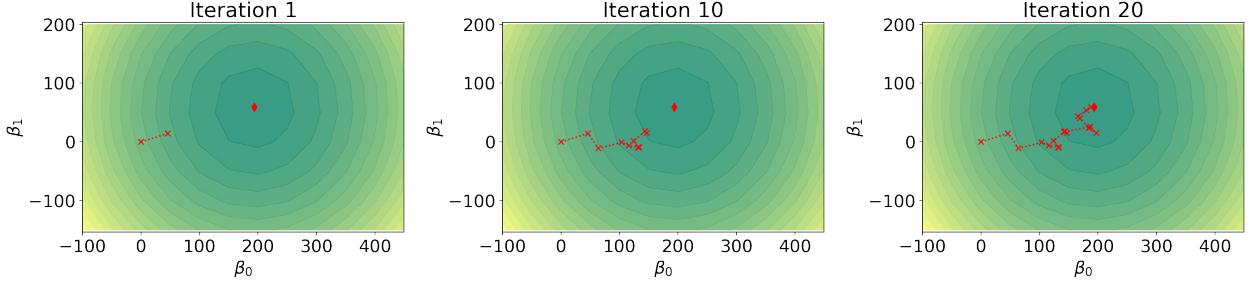


Figure 4.7: Illustration of the stochastic gradient descent algorithm. Contours of the objective function (normalised case), ERM (red diamond), and stochastic gradient descent updates (crosses).

$t$ : iteration of SGD     $i$ :  $i^{\text{th}}$  observation of mini-batch.

*Proof.* Note that, for any  $1 \leq i \leq n_b$ ,  $(\tilde{X}_i^{(t)}, \tilde{Y}_i^{(t)})$  is a uniform discrete random variable taking values in the set  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ . Hence, for any  $1 \leq i \leq n_b$  and  $1 \leq j \leq n$ ,

$$\Pr((\tilde{X}_i^{(t)}, \tilde{Y}_i^{(t)}) = (x_j, y_j)) = \frac{1}{n}.$$

It follows that, for any  $1 \leq i \leq n_b$ ,

$$\mathbb{E}[L(\tilde{Y}_i^{(t)}, h_\theta(\tilde{X}_i^{(t)}))] = \frac{1}{n} \sum_{j=1}^n L(y_j, h_\theta(x_j)) = \hat{R}_n(h_\theta)$$

and

$$\begin{aligned} \mathbb{E}[\nabla_\theta \tilde{J}^{(t)}(\theta)] &= \frac{1}{n_b} \sum_{i=1}^{n_b} \mathbb{E}[\nabla_\theta L(\tilde{Y}_i^{(t)}, h_\theta(\tilde{X}_i^{(t)}))] + \frac{\lambda}{n} \nabla_\theta \text{pen}(h_\theta) \\ &= \frac{1}{n_b} \sum_{i=1}^{n_b} \nabla_\theta \mathbb{E}[L(\tilde{Y}_i^{(t)}, h_\theta(\tilde{X}_i^{(t)})]] + \frac{\lambda}{n} \nabla_\theta \text{pen}(h_\theta) \\ &= \nabla_\theta \hat{R}_n(h_\theta) + \frac{\lambda}{n} \nabla_\theta \text{pen}(h_\theta) = \nabla_\theta J(\theta). \end{aligned}$$

□

For SGD to converge, the learning rate  $\eta_t$  should be such that  $\eta_t \rightarrow 0$ . A standard parametrisation is

$$\eta_t = \frac{\eta_0}{(1 + t/t_0)^a}$$

where  $\eta_0 > 0$ ,  $t_0 > 0$  and  $1/2 < a \leq 1$  are tuning parameters. Alternatively, the learning rate may be tuned adaptively.

**Example 29.** An illustration of stochastic gradient descent, with a batch size of 1, is given in Figure 4.7. After 20 iterations (that is, using only 1/10 of the data), the algorithm has already converged close to the minimum.

## 4.4 Early stopping in (stochastic) gradient descent and implicit regularisation

Consider using (stochastic) gradient descent to minimise the empirical risk  $\hat{R}_n(h_\theta)$  over a large class  $\mathcal{H}$ . As we have seen earlier, without regularisation, the ERM will have a large estimation error, hence a large risk. Early stopping is a strategy that allows to implicitly regularise by stopping the algorithm before it reaches convergence.

Assume we start with an initial prediction rule  $h_{\theta^{(0)}}$  with small complexity (hence large bias). For example, take  $\theta^{(0)} = 0$ . Each step of the (stochastic) gradient descent tends to increase the complexity of the prediction rule, hence to decrease the bias and increase the variance. For small number of iterations  $t$  the bias part dominates, while for large  $t$ , the variance part dominates. One can find the optimal number of steps of the algorithm using a validation set. For each iteration  $t = 0, 1, \dots, t_{\max}$  of the algorithm, we approximate the risk on the validation set, and return the value  $\theta^{(t)}$  that minimises the validation risk.

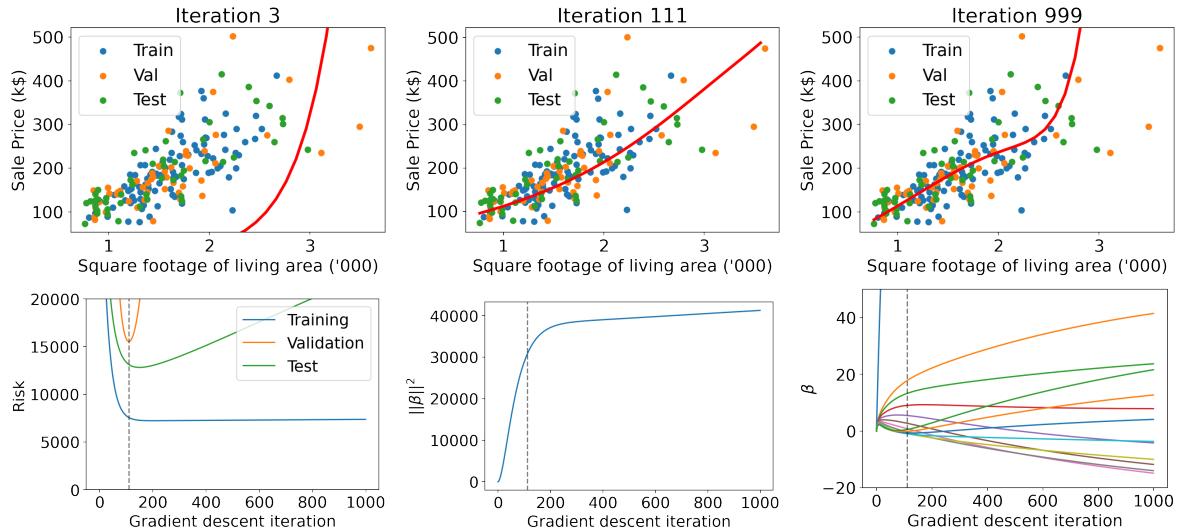


Figure 4.8: Illustration of early stopping and implicit regularisation in gradient descent. Top: Prediction rules at some iterations of the Gradient Descent algorithm. Bottom left: Training, empirical and test risk as a function of the gradient descent iteration. The dotted line corresponds to the iteration where the validation risk is minimised. Bottom middle: Squared norm of beta as a function of the GD iterations. Bottom right: Values of the different dimensions of the parameter  $\beta$  as a function of the number of GD iterations.

**Example 30.** An illustration of early stopping [4.8] is given in Figure [4.8]. On the house pricing dataset, we now fit a polynomial of order 12 (hence  $\beta \in \mathbb{R}^{13}$ ). We split the dataset into a training, validation and test set. We run gradient descent on the training set, starting from the origin. As the number of iterations increases, the norm of the vector  $\beta^{(t)}$  increases. In the first steps, its norm is very small and we are underfitting. For large number of iterations, the iterate  $\beta^{(t)}$  become closer to the minimum, and we start overfitting. The best value is chosen by looking at the empirical risk on the validation set, which is minimised at iteration 111. The learned prediction rule is represented in the top figure (middle).

On large datasets, SGD typically converges faster than GD

Batch size: larger batch sizes lead to a better estimate of the gradient, but are more costly to compute

Randomness of the algorithm allows to escape local minima more easily

Consider using (stochastic) gradient descent to minimise the empirical risk  $\hat{R}_n(h_\theta)$  over a large class  $\mathcal{H}$ .

Without regularisation, the ERM will have a large estimation error, hence a large risk

Early stopping is a strategy that allows to implicitly regularise by stopping the algorithm before it reaches convergence.

# 5 — Linear classifiers

We focus here on binary classification, with  $Y \in \{-1, 1\}$ . Recall that, without loss of generality, a prediction rule can be expressed in terms of a discriminant function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , such that  $h(x) = 1$  if  $f(x) \geq 0$  and -1 otherwise. Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}^p$  be some feature expansion function. Let  $\mathcal{H}$  be the set of classifiers  $h$  with a linear discriminant function

$$\underline{f(x) = \phi(x)^\top \beta}$$

where  $f$  is parameterised by the vector  $\beta \in \mathbb{R}^p$ . Denote  $\mathcal{H}_f = \{f : f(x) = \phi(x)^\top \beta\}$  the set of linear discriminant functions on the transformed inputs  $\phi(x)$ . We have seen in Section 2.2 a generative classifier belonging to this class, the linear discriminant analysis classifier. We describe in this chapter three other linear classifiers.

A natural choice for a classifier would be the empirical risk minimiser under a 0-1 loss with

$$\widehat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i \neq \phi(x_i)^\top \beta}.$$

However, as we discuss in the next section, this classifier is impractical as one cannot usually learn its parameters.

## 5.1 Surrogate loss functions

The 0-1 loss for binary classification is defined as

$$L(y, h(x)) = \mathbb{1}_{y \neq h(x)} = \mathbb{1}_{\text{sign}(yf(x)) = -1} = \psi_{0-1}(yf(x))$$

where  $\psi_{0-1} : \mathbb{R} \rightarrow \{0, 1\}$  is such that  $\underline{\psi_{0-1}(z) = \mathbb{1}_{\text{sign}(z) = -1}}$ . Under the ERM framework, we aim to find the minimiser  $\widehat{f}$  of the empirical risk objective function

$$\frac{1}{n} \sum_{i=1}^n \psi_{0-1}(y_i f(x_i)) \tag{5.1}$$

over some class of functions (e.g. the class of linear functions). The function  $\psi_{0-1}$  is piecewise constant, non-convex, with a discontinuity at 0. The objective function 5.1 is therefore piecewise constant, non-convex and discontinuous, which makes the optimisation problem challenging, and prevents the use of gradient algorithms. To resolve this issue, we typically use an alternative loss function, called a **surrogate loss function**, which will make the objective function easier to minimise.

**Definition 31.** A **surrogate loss function** is a continuous function such that

- 1.  $\psi(x) \geq \psi_{0-1}(x)$  for all  $x$
- 2.  $\psi$  is a convex function.

The function  $\psi$  is also often chosen to be differentiable, although optimisation techniques exist for non-differentiable, convex functions. The ERM under the surrogate loss function  $\psi$  is defined as  $\widehat{h}(x) = 1$  if  $\widehat{f}(x) \geq 0$  and -1 otherwise, where

$$\widehat{f} = \arg \min_{f \in \mathcal{H}_f} \frac{1}{n} \sum_{i=1}^n \psi(y_i f(x_i)).$$

If  $f(x) = \phi(x)^\top \beta$ , this leads to the estimate

$$\widehat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \psi(y_i \phi(x_i)^\top \beta)$$

Here are some standard surrogate loss functions, which are plotted in Figure 5.1.

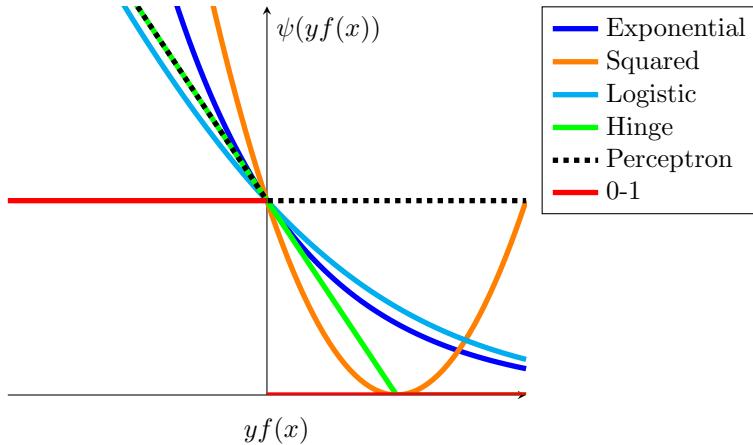


Figure 5.1: Surrogate loss functions for classification.

- Exponential:  $\psi(z) = e^{-z}$
- Squared:  $\psi(z) = (1 - z)^2$
- Logistic:  $\psi(z) = \log(1 + e^{-z}) / \log(2)$
- Perceptron:  $\psi(z) = 1 + \max(0, -z)$
- Hinge:  $\psi(z) = \max(0, 1 - z)$

## 5.2 Least-squares classifier

Consider the squared surrogate loss function

$$\psi(z) = (1 - z)^2. \quad (5.2)$$

The ERM under this surrogate loss and the class of linear classifiers  $\mathcal{H}$  is given by

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (1 - y_i \phi(x_i)^\top \beta)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n y_i^2 (1 - y_i \phi(x_i)^\top \beta)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \phi(x_i)^\top \beta)^2 \end{aligned}$$

$y_i = 1 \quad \theta$ .

=  $\underset{\beta}{\operatorname{argmin}} \frac{1}{n} \| \gamma - \Phi \beta \|^2$   
*Differentiate with respect to  $\beta$ :*  
 $\Phi^\top (\gamma - \Phi \beta) = 0$

and  $\hat{\beta}$  is therefore the typical least-squares estimate

$$\hat{\beta} = (\Phi^\top \Phi)^{-1} \Phi^\top \gamma.$$

This surrogate loss function has the advantage to lead to a closed-form solution. A drawback of the least-squares classifier is that the squared loss is a poor proxy for the 0-1 loss. As is apparent in Figure 5.1, it penalises well-classified points that are far from the boundary, as  $\psi(z)$  takes large values for large positive  $z$ .

**Remark 32.** When the number of examples in each class is the same, the least-squares binary classifier is the same as the binary LDA classifier (see Problem Sheets).

## 5.3 Perceptron

The perceptron algorithm is a classic machine learning classifier algorithm introduced in 1962 by Rosenblatt. Consider the surrogate perceptron loss

$$\underline{\psi(z) = 1 + \max(0, -z)}. \quad (5.3)$$

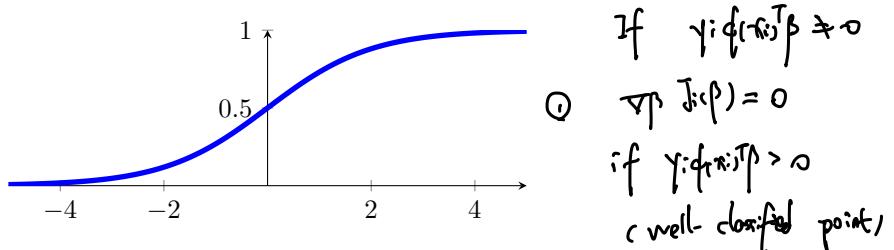


Figure 5.2: Sigmoid function.

This leads to the following objective function

$$J(\beta) = 1 + \frac{1}{n} \sum_{i=1}^n \max(0, -y_i \phi(x_i)^T \beta). \quad (5.4)$$

The minimisation cannot be solved in closed form. The classic perceptron algorithm resorts to a stochastic gradient descent algorithm with a mini-batch size of 1, and a fixed step size  $\eta$  to find the minimiser of (5.4).

### Algorithm 5 Perceptron algorithm

- Initialise  $\beta^{(0)}$  and set  $t = 0$ ;
- For  $t = 1, \dots, n$ 
  - If  $y_t \phi(x_t)^T \beta^{(t)} \geq 0$ , set  $\beta^{(t+1)} = \beta^{(t)}$ ;
  - otherwise, set

$$\beta^{(t+1)} = \beta^{(t)} + \eta y_t \phi(x_t)$$

Note that instead of randomly selecting a batch observation, it just runs over the dataset. The algorithm will converge if the transformed inputs are linearly separable; that is, there exists a function  $h \in \mathcal{H}$  such that  $\hat{R}_n(h) = 0$ .

## 5.4 Logistic regression

Logistic regression is another linear classifier. Like standard linear regression, this classifier can either be interpreted as a plug-in classifier, or as an empirical risk minimiser with a given surrogate function.

**Logistic regression as a plug-in classifier.** Assume that the conditional distribution of  $Y$  given  $X = x$  takes the form

$$\Pr(Y = 1 | X = x) = \text{sig}(\phi(x)^T \beta)$$

and  $\Pr(Y = -1 | X = x) = 1 - \Pr(Y = 1 | X = x) = \text{sig}(-\phi(x)^T \beta)$ , where  $\beta \in \mathbb{R}^p$  is an unknown parameter and

$$\text{sig}(x) = \frac{1}{1 + e^{-x}} \quad (5.5)$$

is the sigmoid function<sup>1</sup> see Figure 5.2. An illustration of the logistic regression model for  $x \in \mathbb{R}^2$  is given in Figure 5.4. Note that we have

$$\log \frac{\Pr(Y = 1 | X = x)}{\Pr(Y = -1 | X = x)} = \phi(x)^T \beta$$

hence the Bayes classifier under this model is linear and can be written as

$$h^*(x) = \begin{cases} 1 & \text{if } f(x) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

where  $f(x) = \phi(x)^T \beta \geq 0$  is the linear discriminant function. The (linear) plug-in classifier is therefore given by

$$\widehat{h}(x) = \begin{cases} 1 & \text{if } \phi(x)^T \widehat{\beta} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

<sup>1</sup>It is also called the logistic function. Not to be confused with the logistic loss function.

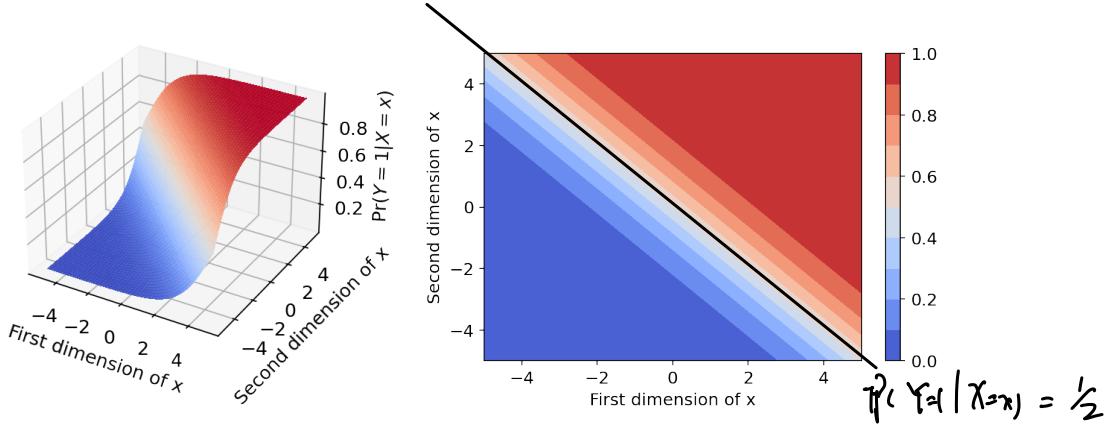


Figure 5.3: Illustration of the logistic regression model, for  $x \in \mathbb{R}^2$ ,  $\phi(x) = x$  and  $\beta = (1 \ 1)^\top$ .

where  $\hat{\beta}$  is some estimate of the parameter  $\beta$ . For instance, one can use a maximum likelihood estimate. The log-likelihood is

$$\begin{aligned}\ell(\beta) &= \sum_{i=1}^n \log \text{sig}(y_i \phi(x_i)^\top \beta) \\ &= - \sum_{i=1}^n \log(1 + e^{-y_i \phi(x_i)^\top \beta})\end{aligned}$$

The MLE is therefore

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^p} - \sum_{i=1}^n \log(1 + e^{-y_i \phi(x_i)^\top \beta}). \quad (5.6)$$

**Logistic regression as a ERM under a surrogate loss function.** The above MLE can be written as

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n \log(2)} \sum_{i=1}^n \log(1 + e^{-y_i \phi(x_i)^\top \beta})$$

hence it can also be interpreted as the ERM under the surrogate logistic loss  $\psi(z) = -\log(\text{sig}(z))/\log(2) = \log(1 + e^{-z})/\log(2)$  amongst the class  $\mathcal{H}$  of linear classifiers.

**Parameter estimation.** Let (dropping the log 2 factor)

$$J(\beta) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \phi(x_i)^\top \beta}) = -\frac{1}{n} \sum_{i=1}^n \log \text{sig}(y_i \phi(x_i)^\top \beta)$$

be the objective function to minimise. First note the following properties of the sigmoid function:

$$\begin{aligned}\text{sig}(-z) &= 1 - \text{sig}(z) \\ \frac{\partial \text{sig}(z)}{dz} &= \text{sig}(z) \text{sig}(-z) \\ \frac{\partial \log \text{sig}(z)}{dz} &= \text{sig}(-z) \\ \frac{\partial^2 \log \text{sig}(z)}{dz^2} &= -\text{sig}(z) \text{sig}(-z)\end{aligned}$$

$J$  is differentiable, with gradient and Hessian

$$\nabla_\beta J(\beta) = -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) \text{sig}(-y_i \phi(x_i)^\top \beta) \quad (5.7)$$

$$\nabla_\beta^2 J(\beta) = \frac{1}{n} \sum_{i=1}^n \text{sig}(y_i \phi(x_i)^\top \beta) \text{sig}(-y_i \phi(x_i)^\top \beta) \phi(x_i) \phi(x_i)^\top \succeq 0 \quad (5.8)$$

We cannot solve  $\nabla_{\beta} J(\beta) = 0$  as there is no closed form solution, and we need to resort to numerical methods. The Hessian matrix  $H(\beta) = \nabla_{\beta}^2 J(\beta)$  is positive semi-definite, and the objective function  $J$  is therefore **convex**. The objective function  $J$  can be minimised using e.g.

- Gradient descent, with update

$$\beta^{(t+1)} = \beta^{(t)} + \frac{\eta}{n} \sum_{i=1}^n y_i \phi(x_i) \text{sig}(-y_i \phi(x_i)^T \beta^{(t)})$$

where  $\eta$  is the learning rate.

- Stochastic gradient descent, with update

$$\beta^{(t+1)} = \beta^{(t)} + \frac{\eta_t}{n_b} \sum_{i=1}^{n_b} \tilde{y}_i^{(t)} \phi(\tilde{x}_i^{(t)}) \text{sig}(-\tilde{y}_i^{(t)} \phi(\tilde{x}_i^{(t)})^T \beta^{(t)})$$

where  $(\tilde{x}_i^{(t)}, \tilde{y}_i^{(t)})$  is the  $t$ th mini-batch of size  $n_b$  and  $(\eta_t)$  is a decreasing sequence of learning rates.

- Newton-Raphson algorithm (known as iterative reweighted least square (IRLS) in this case), with update

$$\beta^{(t+1)} = \beta^{(t)} - (\nabla_{\beta}^2 J(\beta^{(t)}))^{-1} \nabla_{\beta} J(\beta^{(t)}).$$

**Details on the implementation of iterative reweighted least squares.** We can write the gradient and Hessian in a more compact form. Let  $\mu \in \mathbb{R}^n$  be such that  $\mu_i = \text{sig}(\phi(x_i)^T \beta)$  and  $S$  be a  $n$ -by- $n$  diagonal entries with entries  $\mu_i(1 - \mu_i)$ . Let  $c$  be the vector such that  $c_i = \mathbf{1}_{y_i=+1}$ . Then

$$\begin{aligned} n \nabla_{\beta} J(\beta) &= - \sum_{i=1}^n \frac{y_i \phi(x_i) e^{-y_i \phi(x_i)^T \beta}}{1 + e^{-y_i \phi(x_i)^T \beta}} \\ &= \sum_{i=1}^n \phi(x_i)(\mu_i - c_i) \\ &= \Phi^T(\mu - c) \\ n \nabla_{\beta}^2 J(\beta) &= \sum_{i=1}^n \text{sig}(y_i \phi(x_i)^T \beta) \text{sig}(-y_i \phi(x_i)^T \beta) \phi(x_i) \phi(x_i)^T \\ &= \Phi^T S \Phi \end{aligned}$$

Let  $\beta^{(t)}$  be the parameter value after  $t$  iterations of the Newton-Raphson algorithm. Let  $\mu^{(t)}$  and  $S^{(t)}$  be the corresponding vectors and matrices. The Newton-Raphson update is

$$\begin{aligned} \beta^{(t+1)} &= \beta^{(t)} - (\nabla_{\beta}^2 J(\beta^{(t)}))^{-1} \nabla_{\beta} J(\beta^{(t)}) \\ &= \beta^{(t)} + (\Phi^T S^{(t)} \Phi)^{-1} \Phi^T(c - \mu^{(t)}) \\ &= (\Phi^T S^{(t)} \Phi)^{-1} \Phi^T S^{(t)} (\Phi \beta^{(t)} + (S^{(t)})^{-1}(c - \mu^{(t)})) \\ &= (\Phi^T S^{(t)} \Phi)^{-1} \Phi^T S^{(t)} z^{(t)} \end{aligned}$$

where  $z^{(t)} = \Phi \beta^{(t)} + (S^{(t)})^{-1}(c - \mu^{(t)})$ . Then  $\beta^{(t+1)}$  is a solution of the **weighted least squares** problem

$$\begin{aligned} \beta^{(t+1)} &= \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n S_{ii}^{(t)} (z_i^{(t)} - \phi(x_i)^T \beta)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^p} (z^{(t)} - \Phi \beta)^T S^{(t)} (z^{(t)} - \Phi \beta). \end{aligned}$$

Each iteration of the Newton-Raphson algorithm is therefore solving a reweighted least squares problem, hence the name iterative reweighted least square algorithm.

**Linearly separable data.** Assume that the data are linearly separable. That is, there exists a vector  $\tilde{\beta}$  such that  $y_i(\phi(x_i)^T \tilde{\beta}) > 0$  for  $i = 1, \dots, n$ . For any  $c > 0$ , the objective function for  $c\tilde{\beta}$  is

$$J(c\tilde{\beta}) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-cy_i \phi(x_i)^T \tilde{\beta}})$$

which can be made arbitrarily close to 0 as  $c \rightarrow \infty$ . The associated estimated discriminant function  $\hat{f}(x) = \lim_{c \rightarrow \infty} c\phi(x)^T \tilde{\beta}$  is therefore equal to  $+/-\infty$ , depending on the sign of  $\phi(x)^T \tilde{\beta}$ . Similarly, the plug-in estimator of  $\eta(x) = \Pr(Y = 1 | X = x)$  is  $\hat{\eta}(x) = 1$  or 0, yielding overconfident class predictions. One way to address this problem is to add a regularisation term to the objective function.

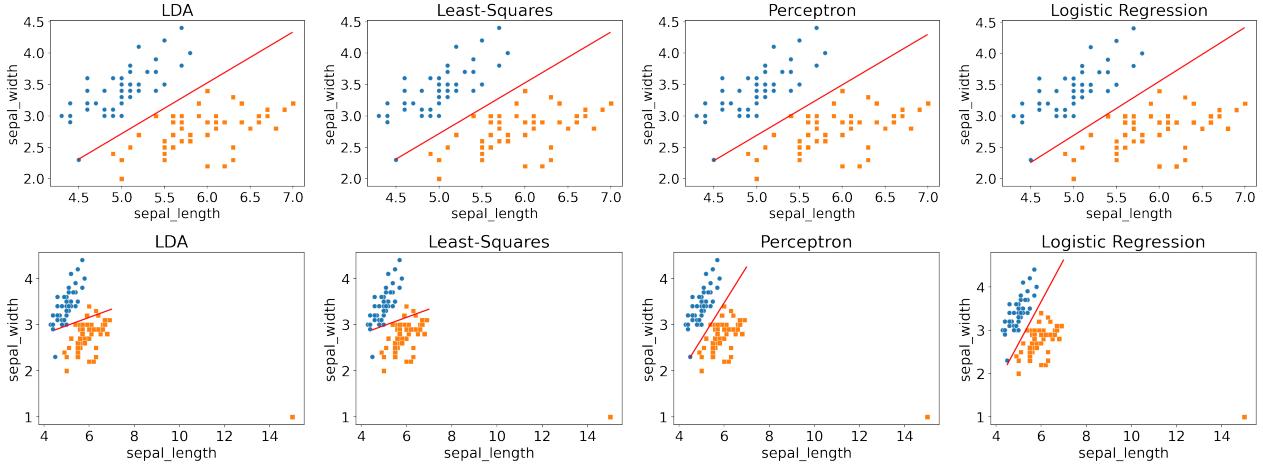


Figure 5.4: Top: Decision boundaries for 4 different linear classifiers on the iris data (2D, sepal length and width), with two classes (setosa and versicolor). Bottom: Same as above, but with the addition of an outlier far from the decision boundary.

**Logistic regression for multi-class classification.** Let  $Y \in \{1, \dots, K\}$ . Multi-class logistic regression uses the **softmax function** to model the conditional class probabilities. For  $k = 1, \dots, K$

$$\Pr(Y = k \mid X = x) = \frac{\exp(\phi(x)^\top \beta_k)}{\sum_{k'=1}^K \exp(\phi(x)^\top \beta_{k'})}$$

where  $\beta_k \in \mathbb{R}^p$  for  $k = 1, \dots, K$ . This yields linear decision boundaries as

$$\log \frac{\Pr(Y = k \mid X = x)}{\Pr(Y = \ell \mid X = x)} = \phi(x)^\top (\beta_k - \beta_\ell).$$

## 5.5 Examples

### 5.5.1 Binary classification on the Iris data

Consider first the iris dataset, studied in Section 2.2. We consider the  $n = 100$  observations from the classes setosa and versicolor, and consider only the sepal length and sepal width. The classification boundaries of linear discriminant analysis, logistic regression, perceptron and least-squares are shown in Figure 5.4 (top). Note that as the classes have the same number of elements, the least-squares classifier and LDA are equivalent. In this example, the data are linearly separable, and all methods give similar results.

To emphasise the differences between the different methods, consider that we have an additional data example from the class versicolor, far from the decision boundary. As shown in Figure 5.4 (bottom), the presence of this outlier has a large influence on the decision boundary for LDA/least-squares. For both the perceptron and logistic regression, the outlier has almost no influence on the decision boundary.

### 5.5.2 Multiclass classification on the Crabs data

In this example, we consider the Crabs dataset ( $K = 4$  classes). We first extract the first two LDA discriminant coordinates of each example, denoted  $z_1, \dots, z_n$ , and use these coordinates as our input data. We first consider a simple linear model, that is a transformation function  $\phi(z_i) = (1 \ z_{i1} \ z_{i2})^\top$ . The learned linear prediction rules for LDA and logistic regression are shown in Figure 5.5 (left-middle). Next, we consider an input expansion that includes interaction terms. That is, we set  $\phi(z_i) = (1 \ z_{i1} \ z_{i2} \ z_{i1}z_{i2})^\top$ . The learned prediction rule, which now has nonlinear decision boundaries in the 2D input space, is represented in Figure 5.5 (right).

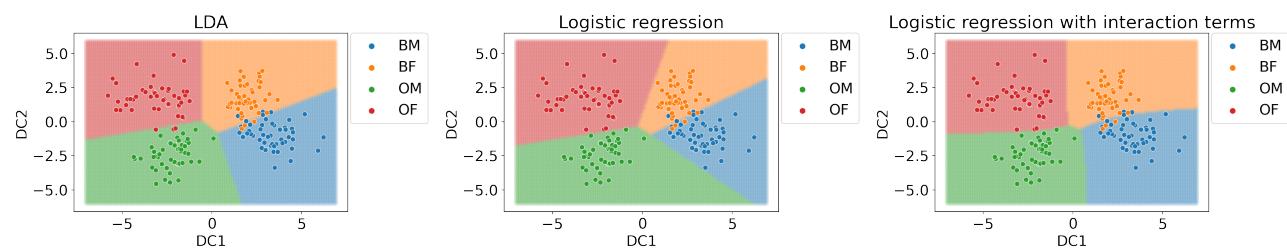


Figure 5.5: Multiclass classification on the Crabs dataset (2D projections), using LDA, logistic regression and logistic regression with interaction terms. The data are represented by dots.

# Statistical Machine Learning

## Hilary Term 2021

§  
3

François Caron  
Department of Statistics  
University of Oxford

Slide credits and other course material can be found at:  
<https://canvas.ox.ac.uk/courses/65441>

## Parametric vs Nonparametric

So far, the methods studied were parametric: prediction rule  $h_\theta$  parameterised by some (finite-dimensional) parameter  $\theta$   
 $\theta$  is fitted using maximum likelihood or empirical risk minimisation using the training data, giving the learned prediction rule  $\hat{h}_\theta$   
In this lecture: Model-free, Nonparametric approach  
Number of parameters grows with the number of data

1

## Nearest neighbor (NN) classification/regression

Let  $\mathcal{X}$  be the input space

Let  $\rho : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$  be some distance/dissimilarity function  
E.g., Euclidian distance (L2 norm)

$$\rho(x, x') = \|x - x'\|_2$$

2

## Nearest neighbor (NN) classification/regression

Training: Store the entire training set  $(x_1, y_1), \dots, (x_n, y_n)$

Learned Prediction rule

$$\hat{h}^{(d)}(x) = y_{\text{nn}(x)}$$

where  $\text{nn}(x) \in \{1, 2, \dots, n\}$ , is the index of the training example whose input is the closest to  $x$

$$\text{nn}(x) = \arg \min_{i \in \{1, 2, \dots, n\}} \|x - x_i\|_2^2$$

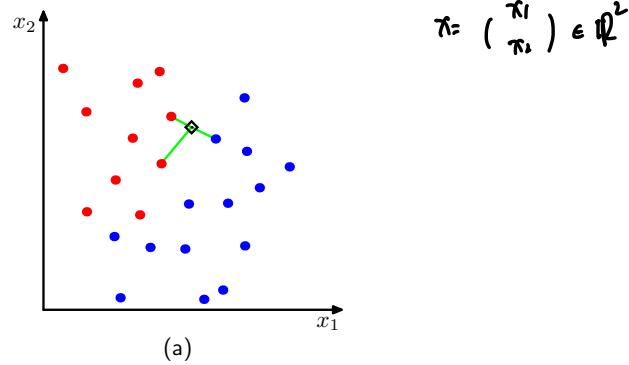
Inductive bias: Label of point is similar to the label of nearby points

3

4

## Visual example

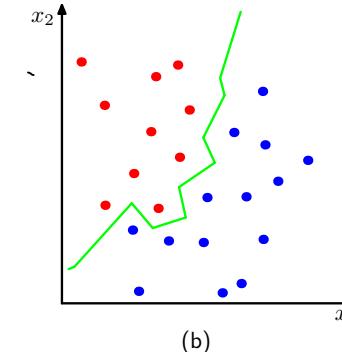
In this 2-dimensional example, the nearest point to  $x$  is a red training instance, thus,  $x$  will be labeled as red.



(a)

## Decision boundary

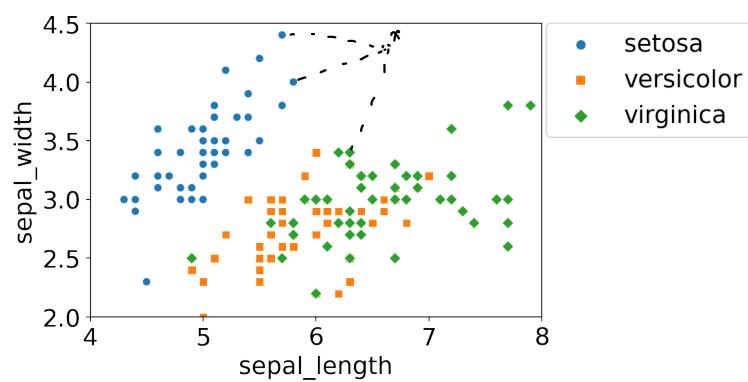
For every point in the space, we can determine its label using the NN classification rule. This gives rise to a decision boundary that partitions the space into different regions.



5

## Nearest neighbor (NN) classification/regression

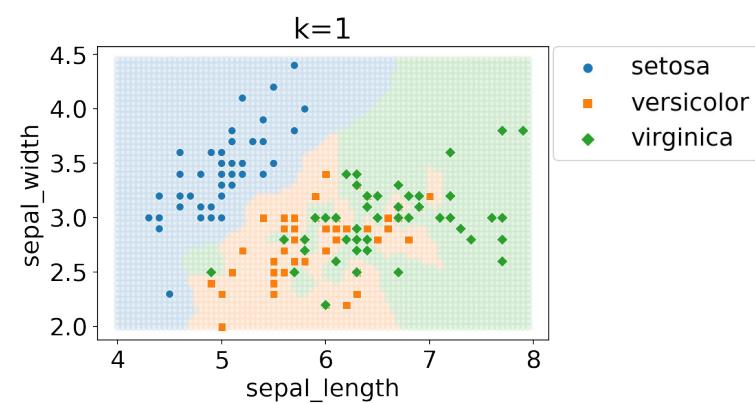
Example on Iris data (2D)



6

## Nearest neighbor (NN) classification/regression

Example on Iris data (2D)



7

7

## Nearest neighbor (NN) classification/regression

### Theorem (Cover-Hart Inequality)

Let  $\hat{h}^{(d_n)}$  be the NN learned classifier for a dataset of size  $n$ . Under mild assumptions,

$$R(h^*) \leq \lim_{n \rightarrow \infty} \mathbb{E}[R(\hat{h}^{(D_n)})] \leq 2R(h^*)(1 - R(h^*)) \leq 2R(h^*)$$

$\leq 1$

The risk is asymptotically at worst twice that of the Bayes optimal classifier.

8

## Regression/classification with $k$ neighbours?

Denote  $\text{knn}(x) \subset \{1, \dots, n\}$  the subset of cardinality  $k$  corresponding to the indices of the  $k$  nearest neighbours of  $x$

Classification Rule ( $y_i \in \{1, \dots, K\}$ )

Every neighbour votes: neighbour  $i \in \text{knn}(x)$  votes for class  $y_i$

Aggregate everyone's vote to obtain the discriminant function: for each class  $c = 1, \dots, K$

$$\widehat{f}_c(x) = \sum_{i \in \text{knn}(x)} \mathbb{1}_{y_i=c}$$

Label with the majority

$$\widehat{h}(x) = \arg \max_{c=1, \dots, K} \widehat{f}_c(x)$$

Regression rule

Average across nearest neighbours:

$$\widehat{h}(x) = \frac{\sum_{i \in \text{knn}(x)} y_i}{k}$$

10

## How to measure nearness with other distances?

Previously, we used the Euclidean distance

$$\text{nn}(x) = \arg \min_{i \in \{1, \dots, n\}} \|x - x_i\|_2^2$$

We can also use alternative distances

$$\|x - x_n\|_q = \left( \sum_j |x_j - x_{nj}|^q \right)^{1/q}$$

for  $q \geq 1$ .

E.g., the  $L_1$  distance for  $q = 1$  (i.e., city block distance, or Manhattan distance)

$$\begin{aligned} \text{nn}(x) &= \arg \min_{i \in \{1, \dots, n\}} \|x - x_i\|_1 \\ &= \arg \min_{i \in \{1, \dots, n\}} \sum_{j=1}^p |x_j - x_{ij}| \end{aligned}$$

9

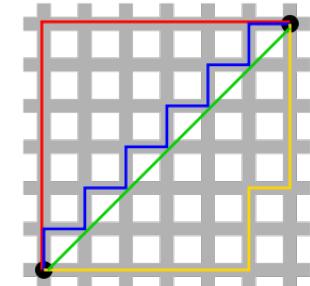


Figure: Green line is Euclidean distance. Red, Blue, and Yellow lines are  $L_1$  distance

## kNN as a plug-in method

Can interpret kNN as a plug-in approach

For classification,  $\Pr(Y = c | X = x)$  is approximated by

$$\frac{1}{k} \sum_{i \in \text{knn}(x)} \mathbb{1}_{y_i=c}$$

For regression, the conditional cdf  $\Pr(Y \leq y | X = x)$  is approximated by

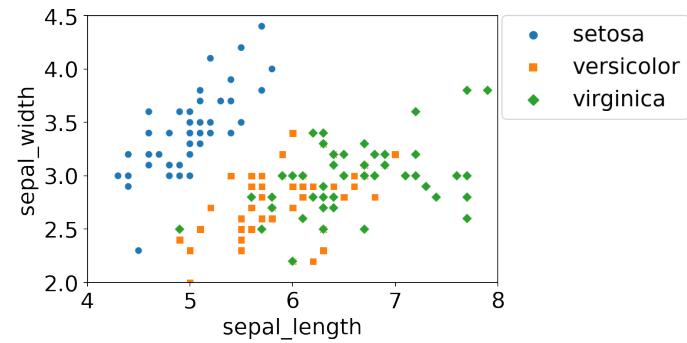
$$\frac{1}{k} \sum_{i \in \text{knn}(x)} \mathbb{1}_{y_i \leq y}$$

Nonparametric estimators of the conditional distribution

11

## Regression/classification with $k$ neighbours?

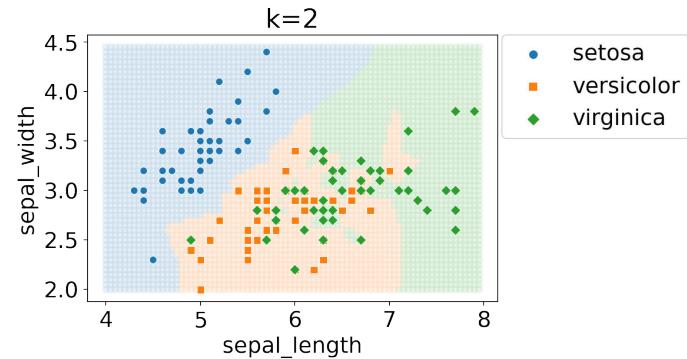
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

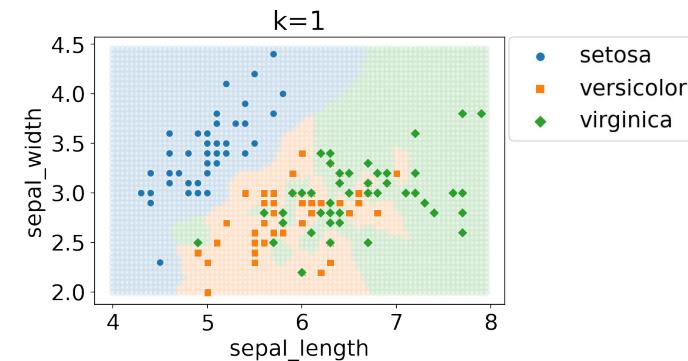
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

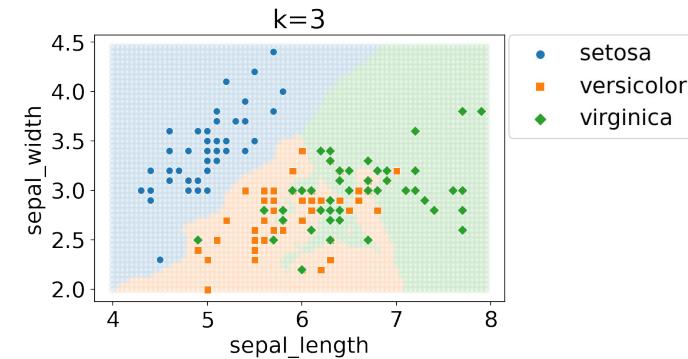
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

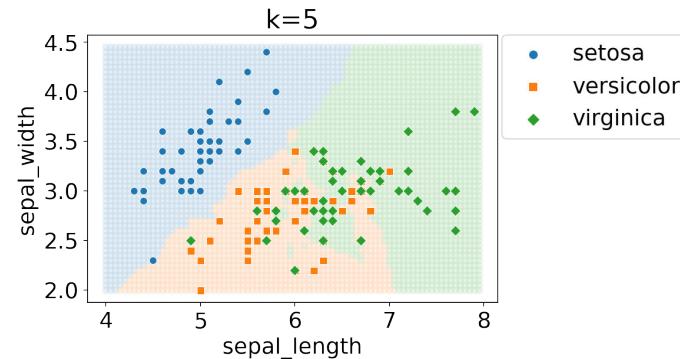
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

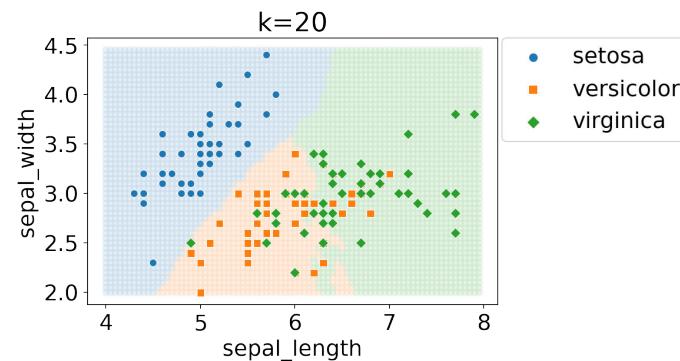
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

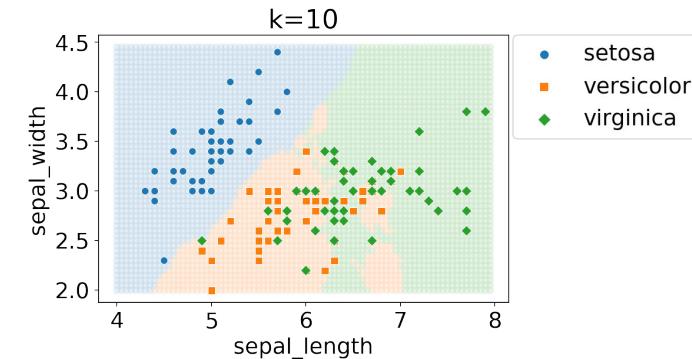
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

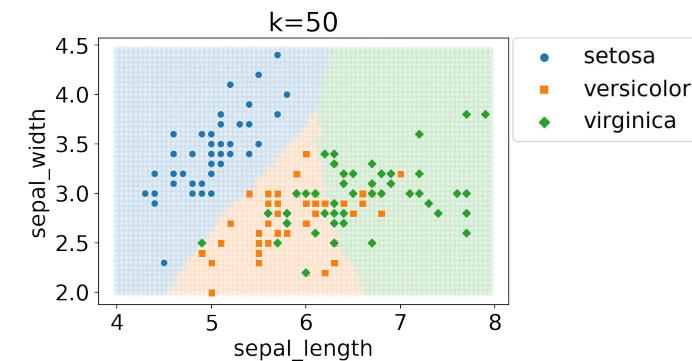
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

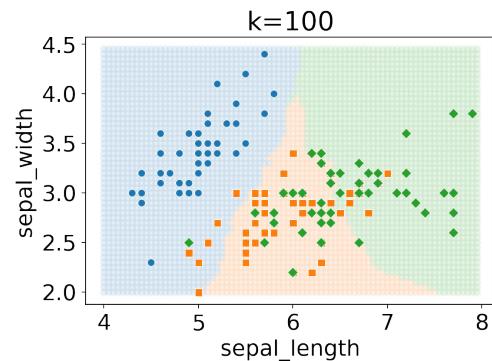
Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

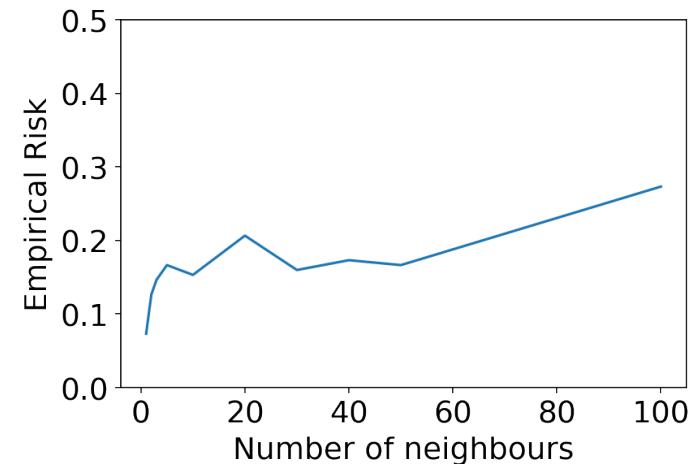
Example on Iris data (2D)



As we increase  $k$ , we get  
a smoother boundary i.e. simpler  
classifier.

## Regression/classification with $k$ neighbours?

Example on Iris data (2D)



12

## Regression/classification with $k$ neighbours?

### Bias/variance - Approximation/estimation tradeoff

Small number of neighbours  $k$ : small bias, large variance

Large number of neighbours  $k$ : large bias, small variance

$k$  may be chosen by using a validation set, or cross-validation

Choosing the right distance is also important

13

## Preprocess data

Assuming all features are equally important!

Distances depend on units of the features.

Normalise data to have zero mean and unit standard deviation in each dimension

Compute the means and standard deviations in each dimension  $j = 1, \dots, p$

$$\bar{x}_j = \frac{1}{n} \sum_i x_{ij}, \quad s_j^2 = \frac{1}{n-1} \sum_i (x_{ij} - \bar{x}_j)^2$$

Scale the feature accordingly

$$x_{ij} \leftarrow \frac{x_{ij} - \bar{x}_j}{s_j}$$

Many other ways of normalising data

14

15

## Summary

### Advantages of k-NN

Simple and easy to implement – just computing distance

Theoretically, has strong guarantees of “doing the right thing”

### Disadvantages of k-NN

Computationally intensive:  $O(np)$  for labeling a data point

Curse of dimensionality: Need a lot of data for large  $p$ . May want to reduce dimensions first.

Not useful for understanding relationships between attributes.

We need to “carry” the training data around (nonparametric approach).

Choosing the right distance measure and  $k$  can be involved.

## Statistical Machine Learning Hilary Term 2021

François Caron

Department of Statistics

University of Oxford

Slide credits and other course material can be found at:

<https://canvas.ox.ac.uk/courses/65441>

16



1



# Neural Networks



Linear prediction rules

Regression

$$h(x) = \phi(x)^T \beta$$

Binary classification

$$h(x) = \begin{cases} 1 & \text{if } \phi(x)^T \beta \geq 0 \\ -1 & \text{Otherwise} \end{cases}$$

Convex optimisation problem to estimate the parameters

Interpretable prediction rule

2

3

## Introduction

Choice of the (fixed) input/feature transformation  $\phi$

Tailored to the problem (images, music, text, etc.)

Difficult and time-consuming to find good features, requires a lot of engineering

Default choice (e.g. polynomial expansions with interactions)

$$\forall i \in \mathbb{N} \quad \phi(x_i) = (1, x_i, x_i^2, x_i^3)$$

$$\vdash x_i = (x_{i1}, x_{i2}) \Leftrightarrow \phi(x_i) = (1, \phi_{i1}, \phi_{i2}, x_{i1} \cdot x_{i2}, \phi_{i1}^2 \cdot \phi_{i2}, \dots)$$

May require a large number of features to represent some sets of functions; e.g. the number of interaction terms of order 3 (of type  $x_{i1}x_{i2}x_{i3}$ ) scales as  $p^3$

4

## Biological inspiration

Basic computational elements:  
neurons.

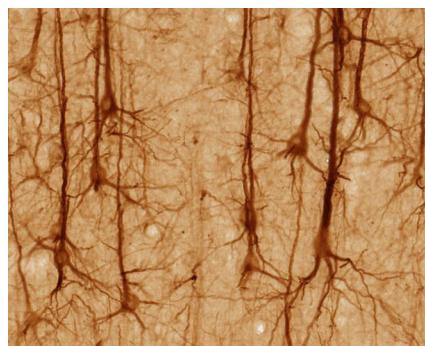
Receives signals from other neurons via dendrites.

Sends processed signals via axons.

Axon-dendrite interactions at synapses.

$10^{10} - 10^{11}$  neurons.

$10^{14} - 10^{15}$  synapses.



6

## Introduction

Alternative is to learn the feature transformations

### Adaptive basis functions

For regression

$$h(x) = \beta_0 + \sum_{k=1}^m \beta_k \phi_{\theta_k}(x)$$

where  $\theta_k$  are unknown parameters.

For classification, discriminant function of the form

$$f(x) = \beta_0 + \sum_{k=1}^m \beta_k \phi_{\theta_k}(x)$$

$\phi_{\theta_k}(x)$  is an adaptive basis function, parameterised by  $\theta_k$ , learned from the data

5

## Neural network for regression

Let  $x_i \in \mathbb{R}^p$  be an input example

Denote  $m \geq 1$  the number of hidden units/neurons

Prediction rule of the form,

$$\text{if } h(x_i) = b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik} \quad \text{output}$$

where for each neuron  $k = 1, \dots, m$

$$z_{ik} = s \left( b_k^{(h)} + \sum_{j=1}^p w_{jk}^{(h)} x_{ij} \right).$$

The nonlinear function  $s$  is known as the activation/transfer function

The vector  $(z_{i1}, \dots, z_{im}) \in \mathbb{R}^m$ , called the vector of hidden units, represents the latent features for the example  $x_i$

potential

7

## Neural network for regression

Parameters that need to be estimated:

Output layer:  $(b^{(o)}, w_1^{(o)}, \dots, w_m^{(o)}) \in \mathbb{R}^{m+1}$

Hidden layer:  $(b_k^{(h)}, w_{1k}^{(h)}, \dots, w_{pk}^{(h)})_{k=1, \dots, m} \in \mathbb{R}^{(p+1)m}$

The intercept terms  $b^{(o)}$  and  $b_k^{(h)}$  are often called the **bias terms** (not to be confused with the bias of an estimator!)

$w_k^{(o)}$  and  $w_{jk}^{(h)}$  are called the **weights**

8

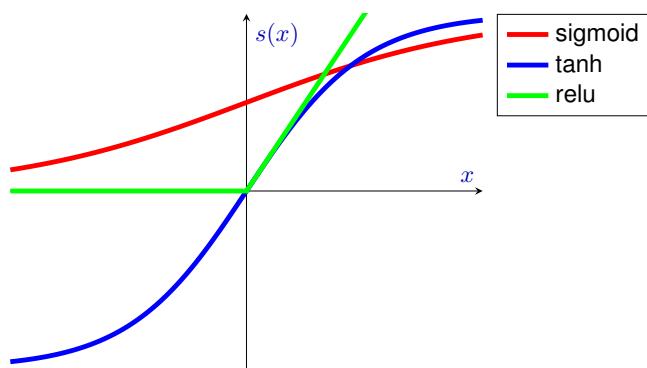
## Neural network for regression

Activation functions

Rectified Linear Unit (ReLU):  $s(x) = \max(0, x)$

sigmoid:  $s(x) = \text{sig}(x) = (1 + e^{-x})^{-1}$

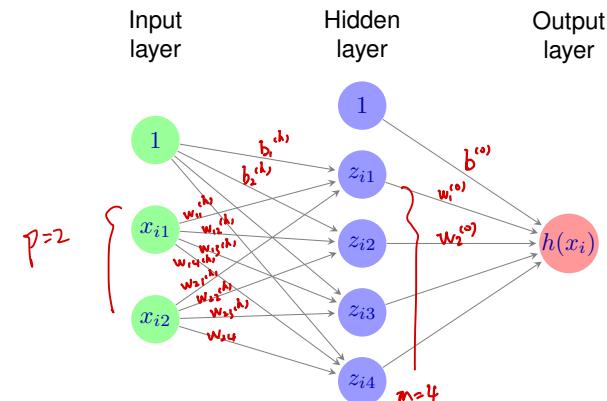
tanh:  $s(x) = \tanh(x)$



10

## Neural network for regression

Illustration with  $p = 2$  and  $m = 4$  hidden inputs/neurons  
4 hidden neurons



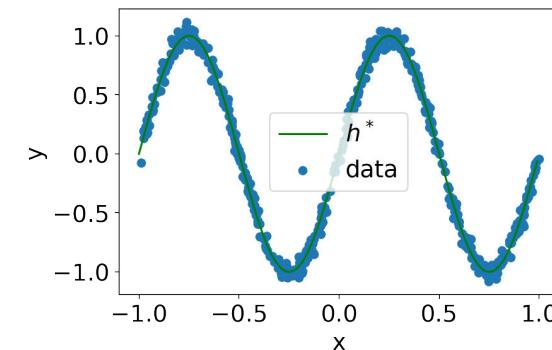
9

## Neural network for regression

Example

Example: synthetic data  $(x_i, y_i) \in \mathbb{R}^2$  ( $n = 500$ ) with Bayes prediction rule

$$h^*(x) = \sin(2\pi x)$$



11

## Neural network for regression

Example

Neural network with  $m = 1, \dots, 5$  hidden units/neurons  
tanh activation function

$$h(x_i) = b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik}$$

where for each neuron  $k = 1, \dots, m$

$$\underline{z_{ik} = \tanh(b_k^{(h)} + w_k^{(h)} x_i) := \phi_k(x_i).}$$

Parameters  $\theta = (b^{(o)}, (w_k^{(o)}, b_k^{(h)}, w_k^{(h)})_{k=1, \dots, m}) \in \mathbb{R}^{3m+1}$   
 ERM under the squared loss

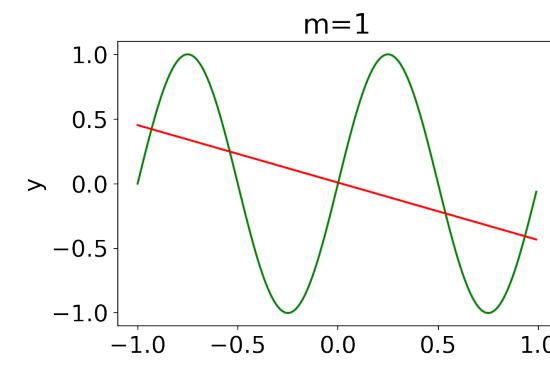
$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - h(x_i))^2$$

(more on the estimation algorithm later)

12

## Neural network for regression

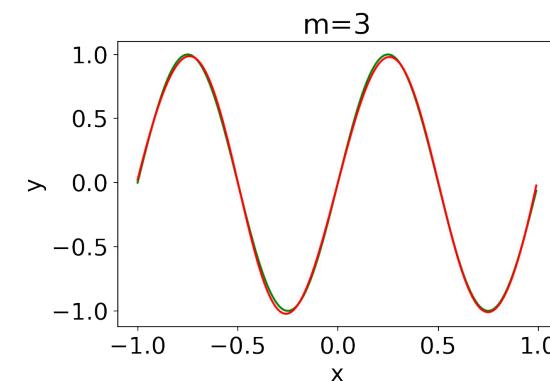
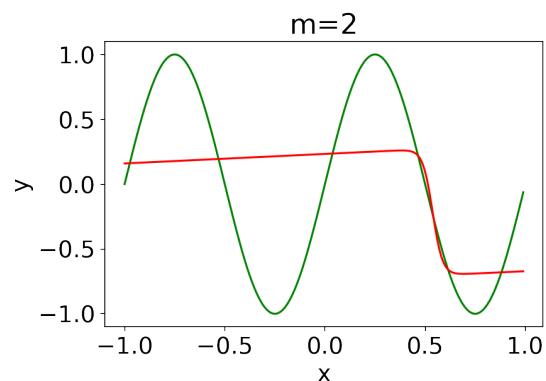
Example



13

## Neural network for regression

Example

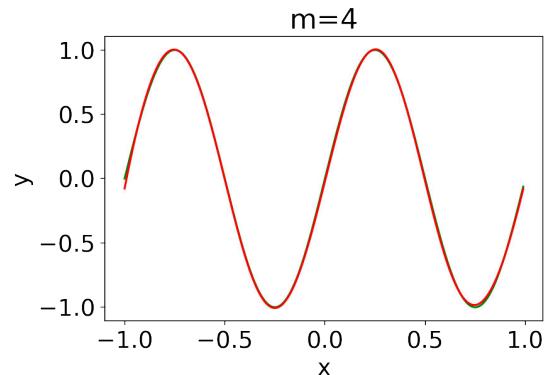


13

13

## Neural network for regression

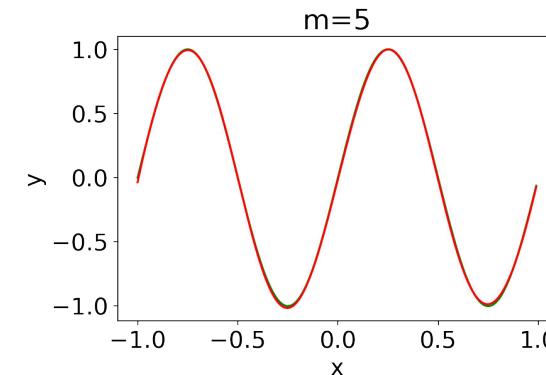
Example



13

## Neural network for regression

Example



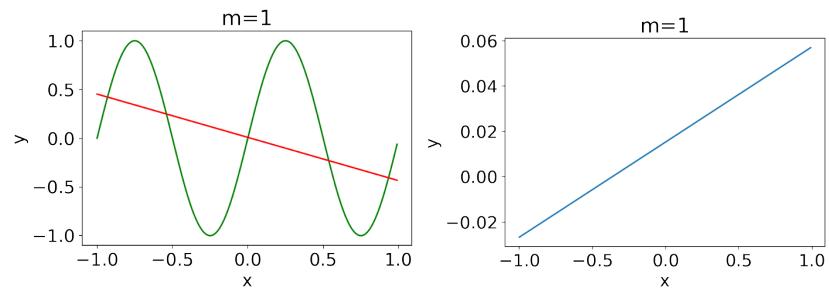
13

## Neural network for regression

Example

Learned input transformations, for  $k = 1, \dots, m$

$$\hat{\phi}_k(x) = \tanh(\hat{b}_k^{(h)} + \hat{w}_k^{(h)}x)$$

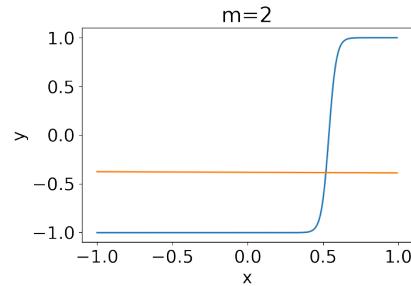
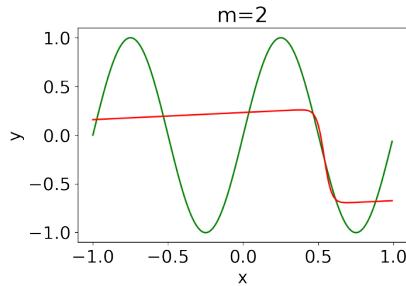


14

15

## Neural network for regression

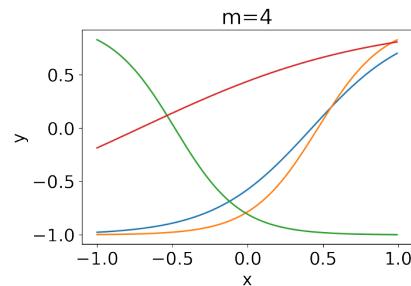
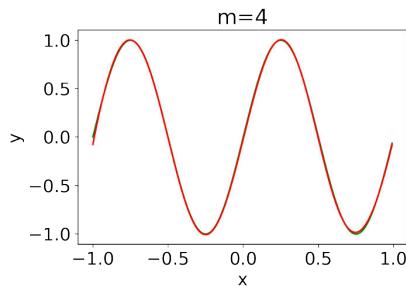
Example



15

## Neural network for regression

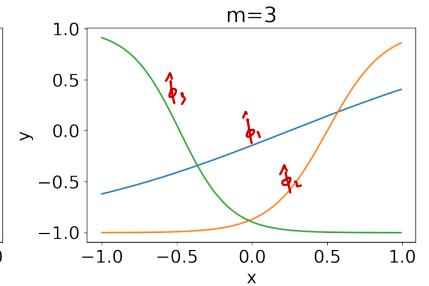
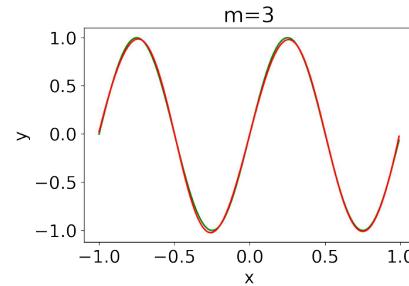
Example



15

## Neural network for regression

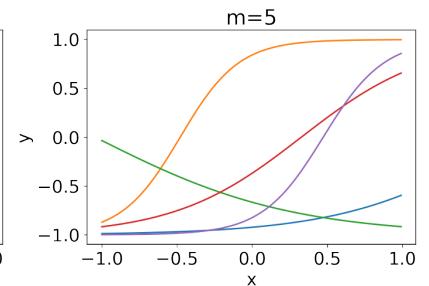
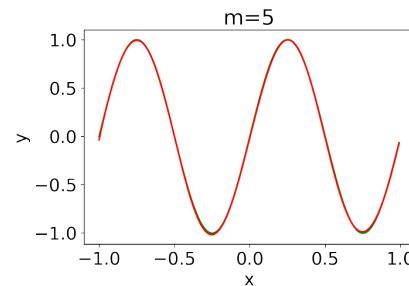
Example



15

## Neural network for regression

Example

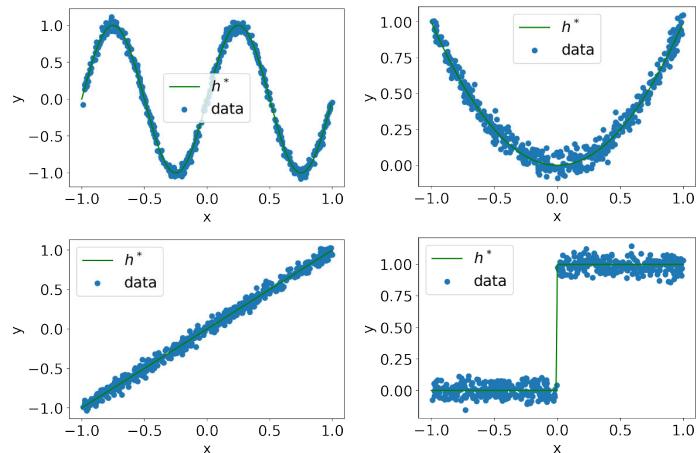


15

## Neural network for regression

Example

Different datasets

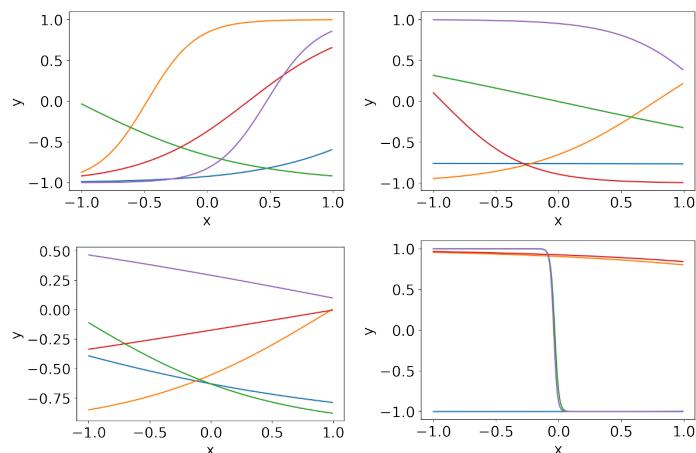


16

## Neural network for regression

Example

Estimated input transformations  $\hat{\phi}_k(x)$ , for  $k = 1, \dots, m$

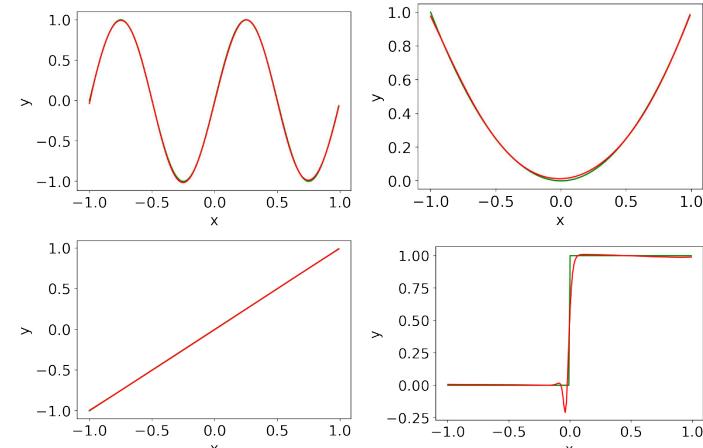


18

## Neural network for regression

Example

Learned prediction rules with  $m = 5$  hidden neurons/units



17

## Neural network for regression

Universal approximation property on  $\mathbb{R}$  (simplified)

Let  $s$  be the relu, tanh or sigmoid activation function. For  $m \geq 1$  let  $\mathcal{F}_{s,m}$  be the set of neural network prediction rules with activation  $s$  and  $m$  hidden units.

**Theorem (Cybenko, 1989)**

Let  $[a, b] \subset \mathbb{R}$  and  $h^*$  be a continuous real-valued function on  $[a, b]$ . For any  $\epsilon > 0$ , there exists  $\underline{m} \geq 1$  and  $h \in \mathcal{F}_{s,\underline{m}}$  such that

$$\underline{|h(x) - h^*(x)| < \epsilon},$$

for all  $x \in [a, b]$ .

By increasing the number  $m$  of neurons in the hidden layer (width) we can approximate any continuous function on a closed interval at an arbitrary precision

Theorem applies more generally to compacts in  $\mathbb{R}^p$ , and activation functions satisfying some mild assumptions

19

## Neural network for binary classification

Let  $x_i \in \mathbb{R}^p$  be an input example

Denote  $m \geq 1$  the number of hidden units/neurons

Prediction rule of the form,

$$h(x_i) = \begin{cases} 1 & \text{If } f(x_i) \geq 0 \\ -1 & \text{Otherwise} \end{cases}$$

with neural network discriminant function

$$f(x_i) = b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik}$$

where for each neuron  $k = 1, \dots, m$

$$z_{ik} = s \left( b_k^{(h)} + \sum_{j=1}^p w_{jk}^{(h)} x_{ij} \right).$$

20

## Neural network for binary classification

Let  $x_i \in \mathbb{R}^p$  be an input example

Denote  $m \geq 1$  the number of hidden units/neurons

Prediction rule of the form,

$$h(x_i) = \begin{cases} 1 & \text{If } \eta(x_i) \geq 1/2 \\ -1 & \text{Otherwise} \end{cases}$$

where

$$\eta(x_i) = \text{sig} \left( b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik} \right)$$

output layer

where for each neuron  $k = 1, \dots, m$

$$z_{ik} = s \left( b_k^{(h)} + \sum_{j=1}^p w_{jk}^{(h)} x_{ij} \right).$$

hidden layer

22

## Neural network for binary classification

Let  $x_i \in \mathbb{R}^p$  be an input example

Denote  $m \geq 1$  the number of hidden units/neurons

Prediction rule of the form,

$$h(x_i) = \begin{cases} 1 & \text{If } \eta(x_i) \geq 1/2 \\ -1 & \text{Otherwise} \end{cases}$$

where  $\eta(x_i) = \text{sig}(f(x_i)) \in [0, 1]$ , with discriminant function

$$f(x_i) = b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik}$$

where for each neuron  $k = 1, \dots, m$

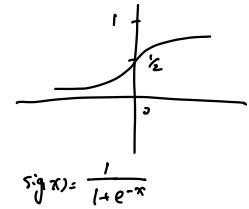
$$z_{ik} = s \left( b_k^{(h)} + \sum_{j=1}^p w_{jk}^{(h)} x_{ij} \right).$$

21

## Neural network for binary classification

Activation function s (relu, tanh, sigmoid) for the hidden layer

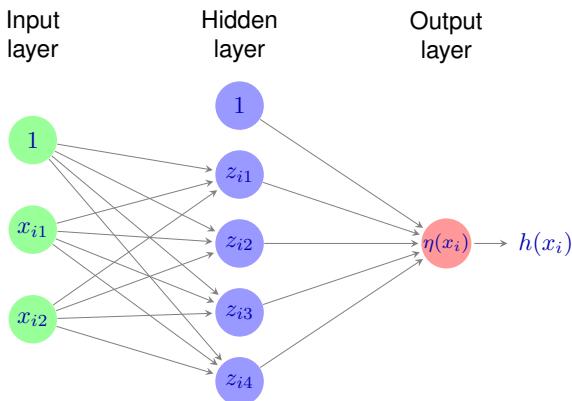
Activation function sig (sigmoid) for the output layer



23

## Neural network for binary classification

Illustration with  $p = 2$  and  $m = 4$  hidden inputs/neurons



24

## Empirical risk minimisation and surrogate loss

Let  $\theta$  be the parameters of the neural network

Under the 0-1 loss, the empirical risk takes the form

$$\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i \neq h_\theta(x_i)}$$

As we have seen earlier, this is not amenable to gradient descent

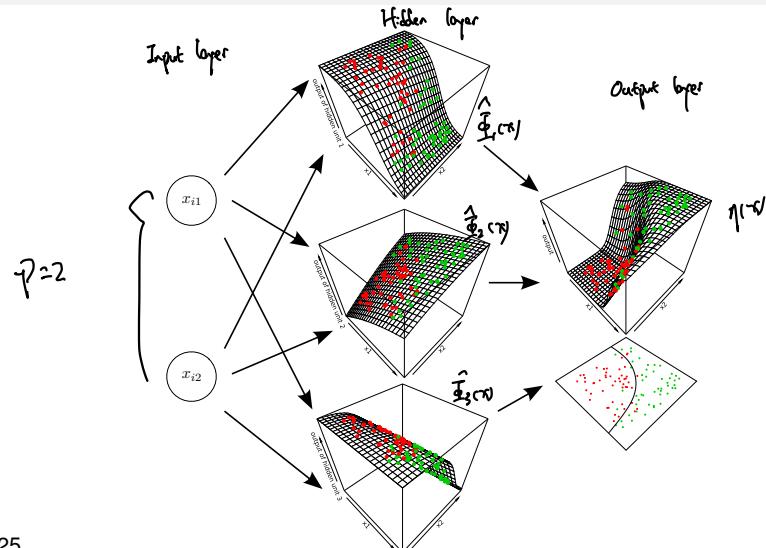
Surrogate logistic loss

$$\eta_\theta(x) = \text{sig}(f_\theta(x))$$

$$\begin{aligned} \psi(y_i f_\theta(x_i)) &= \frac{\log(1 + e^{-y_i f_\theta(x_i)})}{\log(2)} \\ &= \frac{-\log \text{sig}(y_i f_\theta(x_i))}{\log(2)} \\ &= \frac{-\log \eta_\theta(x_i) \mathbb{1}_{y_i=1} - \log(1 - \eta_\theta(x_i)) \mathbb{1}_{y_i=-1}}{\log(2)} \end{aligned}$$

## Neural Network for classification

Illustration for  $p = 2, m = 3$ , with the activation function  $s$  the sigmoid function



25

## Empirical risk minimisation and surrogate loss

Surrogate loss (dropping log 2 factor)

$$\tilde{L}(y_i, \eta_\theta(x_i)) = -\log \eta_\theta(x_i) \mathbb{1}_{y_i=1} - \log(1 - \eta_\theta(x_i)) \mathbb{1}_{y_i=-1}$$

known as the log-loss

ERM under the surrogate loss

$$\begin{aligned} \hat{\theta} &= \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \tilde{L}(y_i, \eta_\theta(x_i)) \\ l_c(\theta) &= -\sum_{i=1}^n \tilde{L}(y_i, \eta_\theta(x_i)) \end{aligned}$$

The learned classifier may equivalently be interpreted as a plug-in classifier under the model

$$\Pr(Y = 1 | X = x) = \eta_\theta(x) = \text{sig}(f_\theta(x))$$

where the parameters  $\theta$  are estimated using maximum-likelihood.

26

27

## Regularisation for neural networks

Parameters

$$\theta = (b^{(o)}, (w_k^{(o)})_{k=1,\dots,m}, (b_k^{(h)})_{k=1,\dots,m}, (w_{jk}^{(h)})_{j=1,\dots,p, k=1,\dots,m}) \in \mathbb{R}^{(p+2)m+1}$$

Prone to over-fitting if  $m$  is large

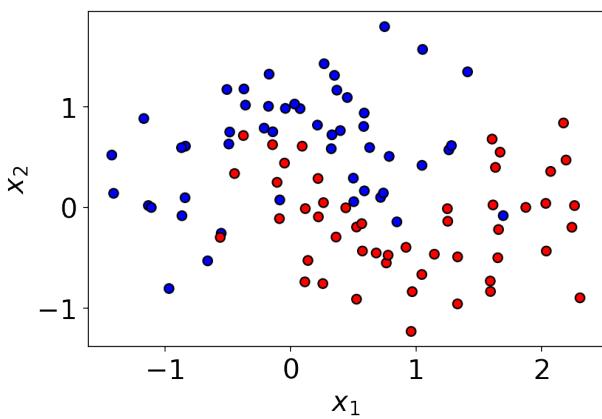
Tikhonov/L2 regularisation on the weights

$$\text{pen}(h_\theta) = \sum_{jk} (w_{jk}^{(h)})^2 + \sum_k (w_k^{(o)})^2$$

Often called weight decay

28

## Illustration: different number of neurons/units



30

## Regularisation for neural networks

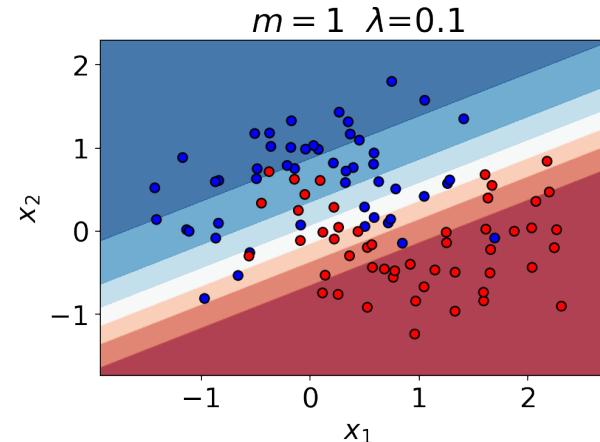
Objective function

$$\begin{aligned} J(\theta) &= \frac{1}{n} \sum_{i=1}^n \tilde{L}(y_i, \eta_\theta(x_i)) + \frac{\lambda}{n} \text{pen}(h_\theta) \\ &= -\frac{1}{n} \sum_{i=1}^n (\log \eta_\theta(x_i) \mathbb{1}_{y_i=1} + \log(1 - \eta_\theta(x_i)) \mathbb{1}_{y_i=-1}) \\ &\quad + \frac{\lambda}{n} \left( \sum_{jk} (w_{jk}^{(h)})^2 + \sum_k (w_k^{(o)})^2 \right) \end{aligned}$$

where  $\lambda \geq 0$  is the regularisation parameter.

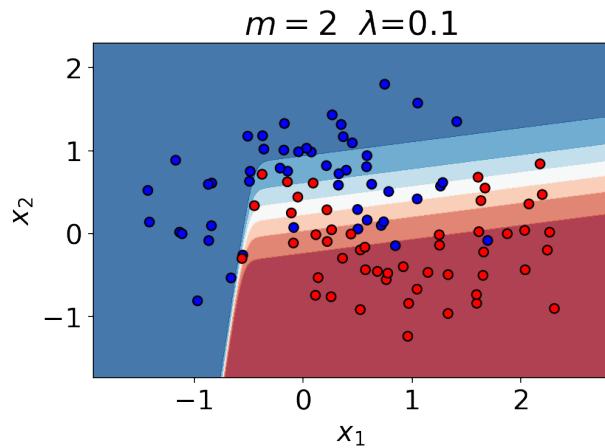
29

## Illustration: different number of neurons/units



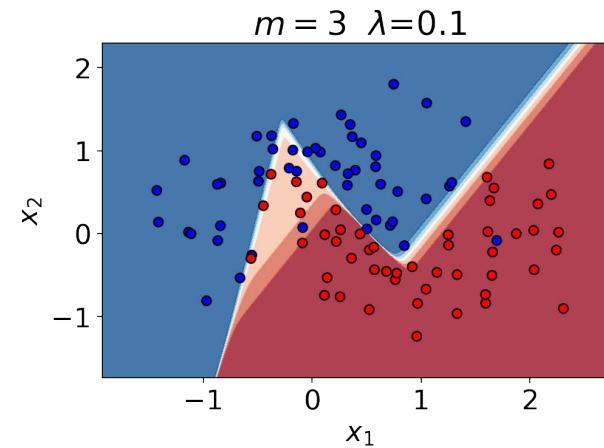
30

## Illustration: different number of neurons/units



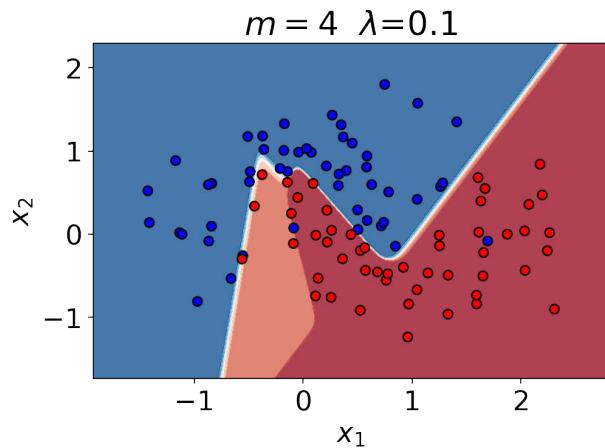
30

## Illustration: different number of neurons/units



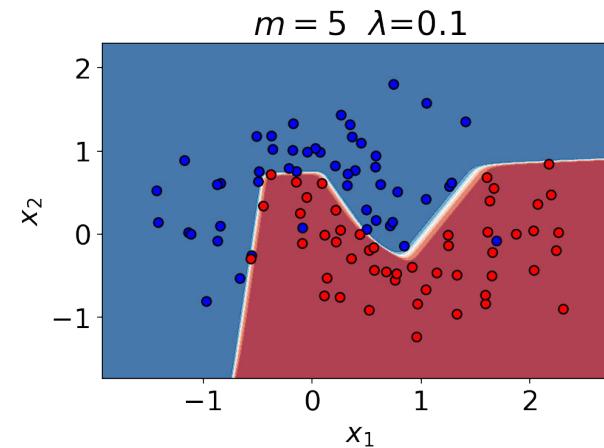
30

## Illustration: different number of neurons/units



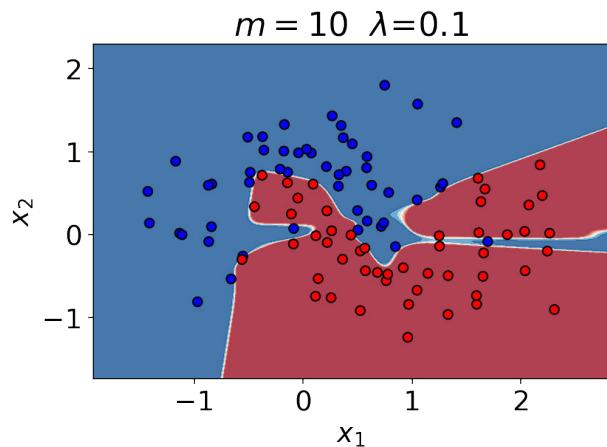
30

## Illustration: different number of neurons/units



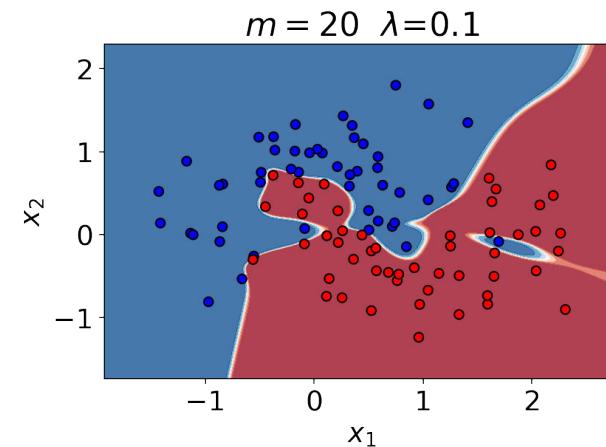
30

## Illustration: different number of neurons/units



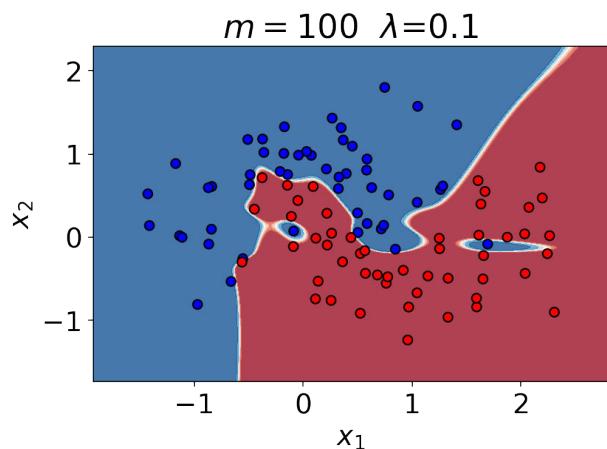
30

## Illustration: different number of neurons/units

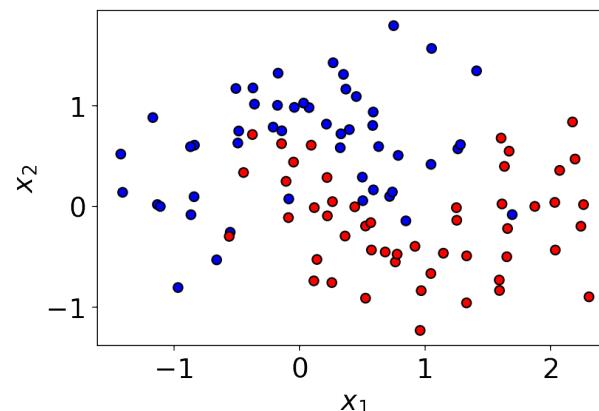


30

## Illustration: different number of neurons/units

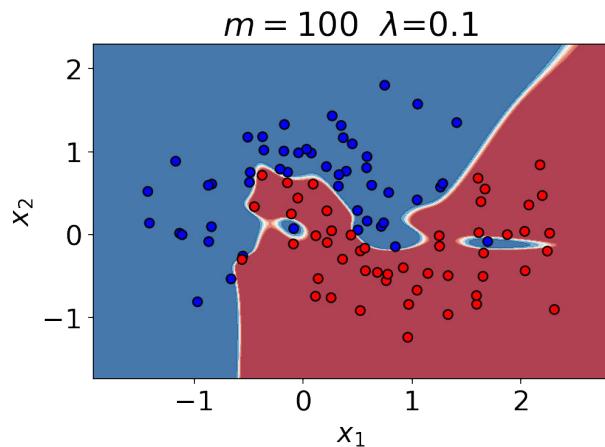


30

Illustration: different values of the regularisation parameter  $\lambda$ 

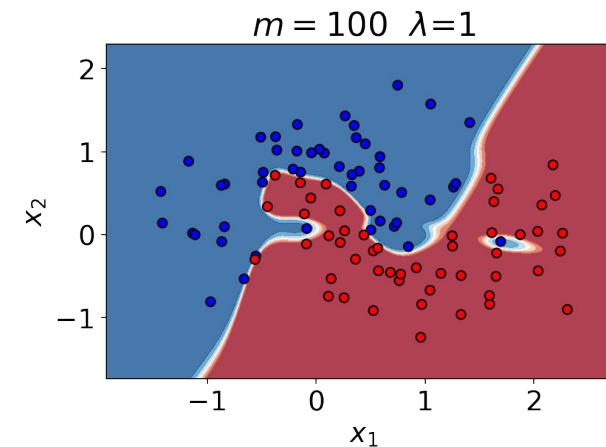
31

Illustration: different values of the regularisation parameter  $\lambda$



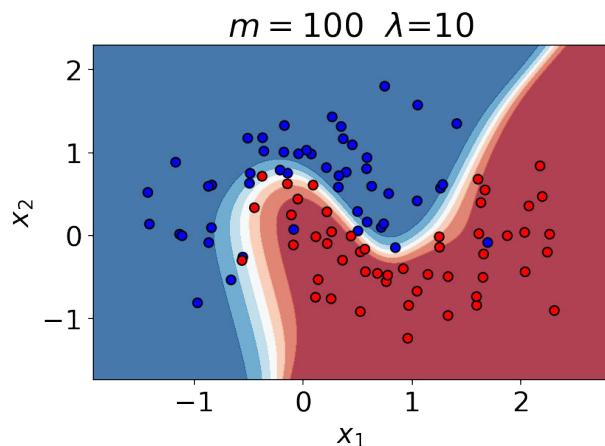
31

Illustration: different values of the regularisation parameter  $\lambda$



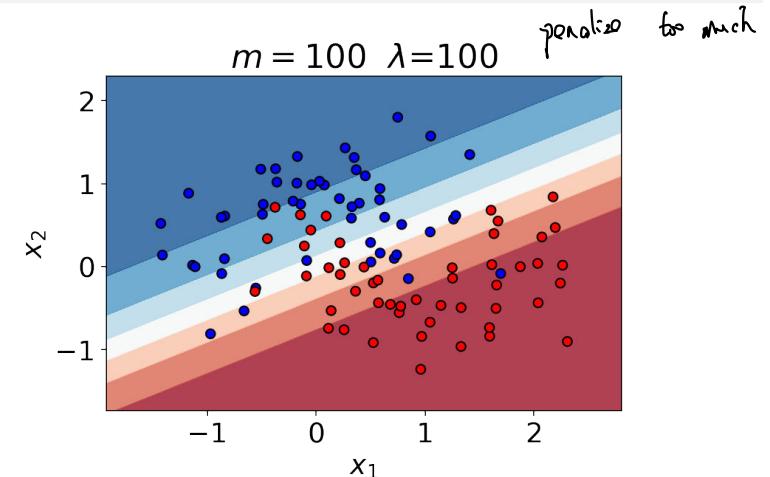
31

Illustration: different values of the regularisation parameter  $\lambda$



31

Illustration: different values of the regularisation parameter  $\lambda$



31

## Training a Neural Network

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(y_i, \eta_\theta(x_i)) + \frac{\lambda}{n} \text{penalty}(\theta)$$

Objective function  $J(\theta)$

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n (\log \eta_\theta(x_i) \mathbb{1}_{y_i=1} + \log(1 - \eta_\theta(x_i)) \mathbb{1}_{y_i=-1}) + \frac{\lambda}{n} \left( \sum_{jk} (w_{jk}^{(h)})^2 + \sum_k (w_k^{(o)})^2 \right)$$


---

Does not admit a closed-form global minimum

Non-convex

Usually resort to (stochastic) gradient descent methods to find a local minimum

32

## Training a Neural Network

Can implement gradient descent or stochastic gradient descent

Due to the non-convexity, different initialisations may lead to different local minima

Early stopping can also be used as implicit regularisation

34

## Training a Neural Network

Assume the activation function of the hidden layer is the sigmoid function  
Using a 0-1 coding  $y_i \in \{0, 1\}$  of the output  
Objective function  $J(\theta)$

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i \log \eta_\theta(x_i) + (1 - y_i) \log(1 - \eta_\theta(x_i))) + \frac{\lambda}{n} \left( \sum_{jk} (w_{jk}^{(h)})^2 + \sum_k (w_k^{(o)})^2 \right)$$

where

$$\eta_\theta(x_i) = \text{sig} \left( b^{(o)} + \sum_{k=1}^m w_k^{(o)} z_{ik} \right) \quad z_{ik} = \text{sig} \left( b_k^{(h)} + \sum_{j=1}^p w_{jk}^{(h)} x_{ij} \right)$$

Using the chain rule (for one variable, see Problem Sheets)

$$\begin{aligned} n \frac{\partial J}{\partial w_k^{(o)}} &= 2\lambda w_k^{(o)} + \sum_{i=1}^n \frac{\partial J}{\partial \eta_\theta(x_i)} \frac{\partial \eta_\theta(x_i)}{\partial w_k^{(o)}} = 2\lambda w_k^{(o)} + \sum_{i=1}^n (\eta_\theta(x_i) - y_i) z_{ik}, \\ n \frac{\partial J}{\partial w_{jk}^{(h)}} &= 2\lambda w_{jk}^{(h)} + \sum_{i=1}^n \frac{\partial J}{\partial \eta_\theta(x_i)} \frac{\partial \eta_\theta(x_i)}{\partial z_{ik}} \frac{\partial z_{ik}}{\partial w_{jk}^{(h)}} \\ &= 2\lambda w_{jk}^{(h)} + \sum_{i=1}^n (\eta_\theta(x_i) - y_i) w_k^{(o)} z_{ik} (1 - z_{ik}) x_{ij}. \end{aligned}$$

33

## Neural Networks for Multiclass classification

$$\mathcal{Y} = \{1, \dots, K\}$$

Use softmax output activations

$$h(x_i) = \arg \max_{c=1, \dots, K} \eta_c(x_i)$$

$$\eta_c(x_i) \in [0, 1] \quad \sum_{c=1}^K \eta_c(x_i) = 1$$

where, for  $c = 1, \dots, K$

$$\eta_c(x_i) = \frac{\exp(b_c^{(o)} + \sum_{k=1}^m w_{kc}^{(o)} z_{ik})}{\sum_{c'} \exp(b_{c'}^{(o)} + \sum_{k=1}^m w_{kc'}^{(o)} z_{ik})}$$

$b_c^{(o)}$  parameters

Surrogate log-loss  $(\eta(x_i) = (\eta_1(x_i), \dots, \eta_K(x_i)))$

at the output layer

$$\tilde{L}(y_i, \eta(x_i)) = - \sum_{c=1}^K \mathbb{1}_{y_i=c} \log \eta_c(x_i)$$

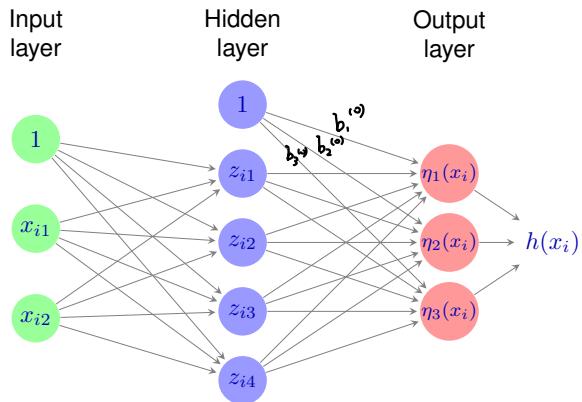
Can be interpreted as maximum likelihood estimation under the model

$$\Pr(Y = c \mid X = x) = \eta_c(x)$$

35

## Neural network for multiclass classification

Illustration with  $p = 2$ ,  $m = 4$  hidden inputs/neurons and  $K = 3$  classes



36

## Shallow neural networks

Neural networks with one hidden layer are called shallow neural networks

Let  $z_i = (z_{i1}, \dots, z_{im})$  be the feature representation of example  $i$

Can be written as (for regression)

$$\underline{h(x_i) = g^{(2)}(z_i), \quad z_i = g^{(1)}(x_i)}$$

where  $g^{(2)}$  and  $g^{(1)}$  correspond to a linear transformation followed by some potentially nonlinear activation

Hence

$$\underline{\underline{h(x_i) = g^{(2)}(g^{(1)}(x_i))}}$$

## Deep learning

37

## Deep neural networks

Deep neural networks are neural networks with more than one hidden layer

Let  $\ell \geq 2$  be the number of layers

$$\underline{\underline{h(x_i) = g^{(\ell+1)}(z_i^{(\ell)})}}$$

for  $\ell = 2, \dots, \ell$

$$\underline{\underline{z_i^{(\ell)} = g^{(\ell)}(z_i^{(\ell-1)})}}$$

and

$$\underline{\underline{z_i^{(1)} = g^{(1)}(x_i)}}$$

Hence

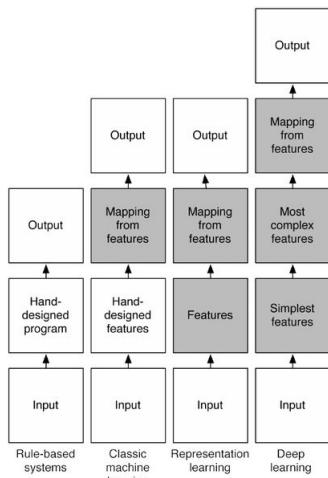
$$\underline{\underline{h(x_i) = g^{(\ell+1)} \left( g^{(\ell)} \left( \dots g^{(1)}(x_i) \right) \right)}}$$

$z_i^{(1)}, \dots, z_i^{(\ell)}$  can be thought as forming a hierarchy from low-level to high-level representations of the example  $x_i$

38

39

## Deep learning intuition



Source: <http://rinuboney.github.io/2015/10/18/theoretical-motivations-deep-learning.html>

40

## Deep neural networks

Why do we need deep networks?

As we have seen, a shallow network ( $\ell = 1$ ) with a large number of neurons  $m$  can approximate any continuous function with compact support

Two issues

$m$  may need to be very large

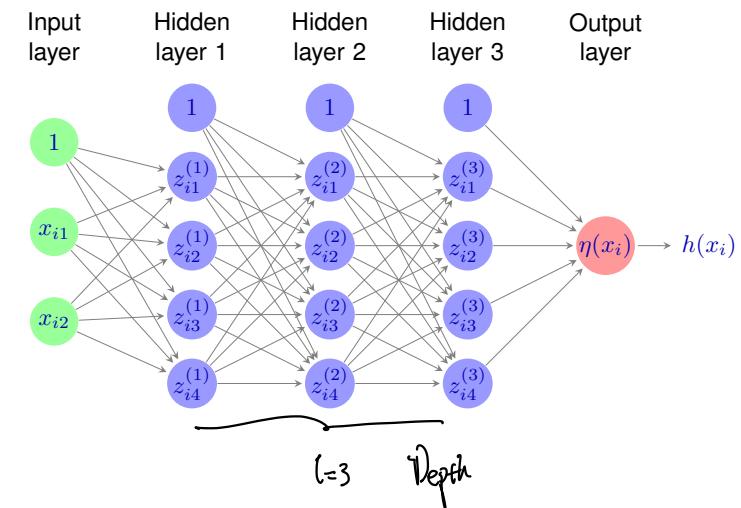
We may not be able to estimate the parameters

Deep networks often lead to more parsimonious representations, and models whose parameters are easier to learn

41

## Feedforward neural network

Illustration with  $p = 2$ ,  $\ell = 3$  layers and  $m = 4$  hidden inputs/neurons per layer



43

42

## Feedforward neural network

A feedforward neural network is a neural network where each layer is fully connected to the layer above/below

Sometimes called multilayer perceptron, but we will not used this name

Let  $\ell$  be the number of hidden layers, and  $m^{(l)}$  be the number of hidden neurons/units at layer  $l = 1, \dots, \ell$

For binary classification

$$\eta(x_i) = \text{sig} \left( b^{(\ell+1)} + \sum_{k=1}^{m^{(\ell)}} w_k^{(\ell+1)} z_{ik}^{(\ell)} \right)$$

and for  $l = 2, \dots, \ell$ ,  $k = 1, \dots, m^{(l)}$

$$z_{ik}^{(l)} = s \left( b_k^{(l)} + \sum_{k'=1}^{m^{(l-1)}} w_{k'k}^{(l)} z_{ik'}^{(l-1)} \right)$$

and for  $k = 1, \dots, m^{(1)}$

$$z_{ik}^{(1)} = s \left( b_k^{(1)} + \sum_{j=1}^p w_{jk}^{(1)} x_{ij} \right)$$

## Feedforward neural network

Matrix notation

For  $l = 1, \dots, \ell$ , define  $\underline{z}_i^{(l)} = (z_{i1}^{(l)}, \dots, z_{im^{(l)}}^{(l)})^\top \in \mathbb{R}^{m^{(l)}}$ , and for  $l = 1, \dots, \ell + 1$  let  $\underline{b}^{(l)} = (b_1^{(l)}, \dots, b_{m^{(l)}}^{(l)})$  and  $\underline{W}^{(l)}$  the  $m^{(l-1)}$ -by- $m^{(l)}$  matrix with entries  $w_{kk'}^{(l)}$ , with  $m^{(\ell+1)} = 1$  and  $m^{(0)} = p$

Denote  $\underline{s}$  the multivariate function where the nonlinear activation  $s$  is applied elementwise

Then

$$\eta(x_i) = \text{sig}\left(b^{(\ell+1)} + (\underline{W}^{(\ell+1)})^\top \underline{z}_i^{(\ell)}\right)$$

and for  $l = 2, \dots, \ell$

$$\underline{z}_i^{(l)} = \underline{s}\left(b^{(l)} + (\underline{W}^{(l)})^\top \underline{z}_i^{(l-1)}\right)$$

and

$$\underline{z}_i^{(1)} = \underline{s}\left(b^{(1)} + (\underline{W}^{(1)})^\top x_i\right)$$

44

## Recap: Chain rule

Let

$$F(x, y) = f(u_1(x, y), \dots, u_p(x, y)) = f(z_1, \dots, z_p)$$

where  $z_j = u_j(x, y)$  for  $j = 1, \dots, p$ .

Chain rule

$$\frac{\partial F}{\partial x} = \sum_{j=1}^p \frac{\partial F}{\partial z_j} \frac{\partial z_j}{\partial x}$$

$$\frac{\partial F}{\partial y} = \sum_{j=1}^p \frac{\partial F}{\partial z_j} \frac{\partial z_j}{\partial y}$$

46

## Learning in feedforward neural networks

Objective function (ignoring the regularisation term for simplicity)

$$\cancel{J(\theta) = \frac{1}{n} \sum_{i=1}^n -\log(\eta_\theta(x_i)) \mathbb{1}_{y_i=1} - \log(1 - \eta_\theta(x_i)) \mathbb{1}_{y_i=-1}}$$

where  $\theta$  is the set of parameters, consisting of the intercepts/biases terms and weights at each layer

No closed-form solution

Nonconvex: Objective functions can have many local minima

On large scale problems, usually use stochastic gradient descent, along with a whole host of techniques for optimization, regularization, and initialization.

Efficient computation of the (stochastic) gradient using backpropagation

45

## Backpropagation

Using 0-1 coding  $y_i \in \{0, 1\}$

Set  $\underline{z}_i^{(\ell+1)} = \eta_\theta(x_i)$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n -y_i \log(z_i^{(\ell+1)}) - (1 - y_i) \log(1 - z_i^{(\ell+1)})$$

Gradients wrt  $z_{ik}^{(l)}$  computed by recursive applications of chain rule, and propagated through the network backwards.

$$\begin{aligned} \frac{\partial J}{\partial z_i^{(\ell+1)}} &= -\frac{1}{n} \left( \frac{y_i}{z_i^{(\ell+1)}} - \frac{(1 - y_i)}{1 - z_i^{(\ell+1)}} \right) \\ \frac{\partial J}{\partial z_{ik}^{(l)}} &= \sum_{r=1}^{m^{(l+1)}} \frac{\partial J}{\partial z_{ir}^{(l+1)}} \frac{\partial z_{ir}^{(l+1)}}{\partial z_{ik}^{(l)}} \\ \frac{\partial J}{\partial w_{jk}^{(l)}} &= \sum_{i=1}^n \frac{\partial J}{\partial z_{ik}^{(l)}} \frac{\partial z_{ik}^{(l)}}{\partial w_{jk}^{(l)}} \end{aligned}$$

47

## Deep learning demo

<http://playground.tensorflow.org/>

48

Deep learning Convolutional neural networks

## Convolutional neural networks

Example: Imagenet<sup>1</sup>

Image classification task

1000 classes, 1.3M training, 100k test



<sup>1</sup>[Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. IJCV. 2015]

## Convolutional neural networks

Convolutional neural networks (CNN) are a specialised type of neural network for dealing with dataset with a grid-like representation

Applications to images, audio, video

State-of-the-art results in the last ten years

49

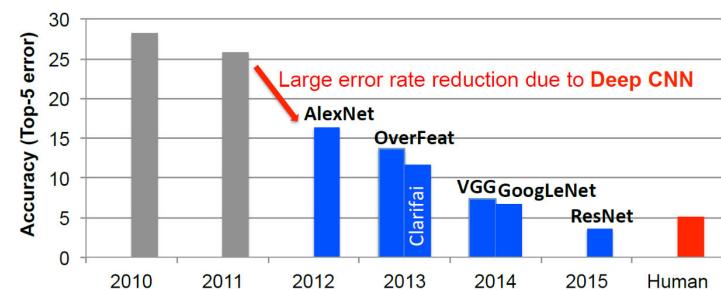
Deep learning Convolutional neural networks

## Convolutional neural networks

Example: Imagenet<sup>1</sup>

Image classification task

1000 classes, 1.3M training, 100k test



<sup>1</sup>[Russakovsky et al. ImageNet Large Scale Visual Recognition Challenge. IJCV. 2015]

## Convolutions

The diagram illustrates the convolution operation between a 7x7 input image and a 3x3 filter. The input image is shown as a grid of 49 pixels, with a 3x3 kernel highlighted in red. The filter is shown as a 3x3 matrix. The result of the convolution is a 3x3 feature map, where each element is the sum of the products of the corresponding kernel element and the input image elements. The resulting feature map is shown as a 3x3 grid.

## Convolutions

$$\left( \begin{array}{cccccc}
 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 x^1 & x^0 & x^1 & x^1 & & & \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 x^0 & x^1 & x^1 & x^0 & & & \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 x^1 & x^0 & x^0 & x^1 & & & \\
 0 & 0 & 0 & 1 & 1 & 0 & 0 \\
 0 & 0 & -0.1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 1 & 1 & 0 & 0 & 0 \\
 0 & 1 & 1 & 0 & 0 & 0 & 0 \\
 1 & 1 & 0 & 0 & 0 & 0 & 0
 \end{array} \right) * \left( \begin{array}{ccc}
 1 & 0 & 1 \\
 0 & 1 & 0 \\
 1 & 0 & 1
 \end{array} \right) = \left( \begin{array}{ccccc}
 1 & 4 & 3 & 4 & 1 \\
 1 & 2 & 4 & 3 & 3 \\
 1 & 2 & 3 & 4 & 1 \\
 1 & 3 & 3 & 1 & 1 \\
 3 & 3 & 1 & 1 & 0
 \end{array} \right)$$

51

## Convolutions

The diagram illustrates the convolution process. A 7x7 input image (labeled "7-by-7 Image") is multiplied by a 3x3 filter (labeled "3-by-3 Filter"). The result is a 3x3 feature transformation map (labeled "Feature transformation/map").

The input image is shown as a grid of values. A 3x3 subgrid in the top-left corner is highlighted with a red border, representing the receptive field of the top-left unit in the output map. The filter is shown as a 3x3 matrix with values 1, 0, 1; 0, 1, 0; 1, 0, 1. The resulting output map shows values 1, 4, 3; 1, 2, 4; 1, 2, 3.

## Convolutions

$$\left( \begin{array}{ccccccc}
 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 1 & 1 & 0 & 0 & 0 \\
 1 & 1 & 0 & 0 & 0 & 0 & 0
 \end{array} \right) * \left( \begin{array}{ccc}
 1 & 0 & 1 \\
 0 & 1 & 0 \\
 1 & 0 & 1
 \end{array} \right) = \left( \begin{array}{ccccc}
 1 & 4 & 3 & 4 & 1 \\
 1 & 2 & 4 & 3 & 3 \\
 1 & 2 & 3 & 4 & 1 \\
 1 & 3 & 3 & 1 & 1 \\
 3 & 3 & 1 & 1 & 0
 \end{array} \right)$$

51

## Convolutions

$$\begin{pmatrix}
 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 0 & 0 & 0 & 1 & 1 & 0 & 0 \\
 0 & 0 & 1 & 1 & 0 & 0 & 0 \\
 0 & 1 & 1 & 0 & 0 & 0 & 0 \\
 1 & 1 & 0 & 0 & 0 & 0 & 0
 \end{pmatrix} \times \begin{pmatrix}
 1 & 0 & 1 \\
 0 & 1 & 0 \\
 1 & 0 & 1
 \end{pmatrix} = \begin{pmatrix}
 1 & 4 & 3 & 4 & 1 \\
 1 & 2 & 4 & 3 & 3 \\
 1 & 2 & 3 & 4 & 1 \\
 1 & 3 & 3 & 1 & 1 \\
 3 & 3 & 1 & 1 & 0
 \end{pmatrix}$$

7-by-7 Image      3-by-3 Filter      Feature transformation/map

## Convolutions

The diagram illustrates the convolution operation. On the left, a **7-by-7 Image** is shown as a matrix of 49 elements. A red box highlights the top-left 3x3 submatrix:  $\begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$ . Above this submatrix, the labels  $x_1 \times x_0 \times x_1$  are written. To its right, the labels  $1 \times 1 \times 1$  are written. Below the submatrix, the labels  $x_0 \times x_1 \times x_0$  are written. The entire 7x7 matrix has labels  $x_1 \times x_0 \times x_1$  written along its bottom and right edges. In the center, a **\*** symbol indicates multiplication. To the right of the **\*** symbol is a **3-by-3 Filter** represented as a 3x3 matrix:  $\begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$ . To its right, the label **=** is followed by a **Feature transformation/map** represented as a 3x3 matrix:  $\begin{pmatrix} 1 & 4 & 3 & 4 & 1 \\ 1 & 2 & 4 & 3 & 3 \\ 1 & 2 & 3 & 4 & 1 \\ 1 & 3 & 3 & 1 & 1 \\ 3 & 3 & 1 & 1 & 0 \end{pmatrix}$ . In the resulting matrix, the second row, third column element (value 1) is highlighted with a green box.

51

## Convolutions

$$\left( \begin{array}{cccccc}
 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 \\
 0 & 0 & 0 & 1 & x_0 & x_0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 0 & 0 \\
 0 & 0 & 1 & 1 & 0 & 0 & 0 \\
 0 & 1 & 1 & 0 & 0 & 0 & 0 \\
 1 & 1 & 0 & 0 & 0 & 0 & 0
 \end{array} \right) * \left( \begin{array}{ccc}
 1 & 0 & 1 \\
 0 & 1 & 0 \\
 1 & 0 & 1
 \end{array} \right) = \left( \begin{array}{ccccc}
 1 & 4 & 3 & 4 & 1 \\
 1 & 2 & 4 & 3 & 3 \\
 1 & 2 & 3 & 4 & 1 \\
 1 & 3 & 3 & 1 & 1 \\
 3 & 3 & 1 & 1 & 0
 \end{array} \right)$$

7-by-7 Image      3-by-3 Filter      Feature transformation/map

## Convolutions

$$\left( \begin{array}{ccccccc} 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \end{array} \right) * \left( \begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccccc} 1 & 4 & 3 & 4 & 1 \\ 1 & 2 & 4 & 3 & 3 \\ 1 & 2 & 3 & 4 & 1 \\ 1 & 3 & 3 & 1 & 1 \\ 3 & 3 & 1 & 1 & 0 \end{array} \right)$$

51

## Convolutions

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

Click for animation.

Source: <https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/>

52

## Summary

Deep neural networks are heavily parameterised models that learn hierarchical representations of the data

Model parameters are learned using stochastic gradient descent

Regularisation via L2 penalisation (weight decay), early stopping, or others (e.g. dropout, not covered here)

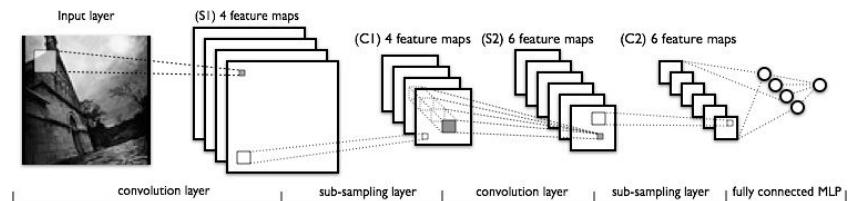
Architecture of the network is important: type of layers (fully connected, convolutional, etc.), number of hidden layers, number of neurons per layers; can be chosen using a validation set

Requires a lot of data and computing resources

State-of-the-art results for many applications in computer vision, natural language processing and others

Python libraries: pytorch, tensorflow, keras

## Deep Convolutional Neural Networks



Input: 2D image.

### Convolution layer

Applies different filters, whose weights are learned from the data, followed by a nonlinear activation

Can be seen as a neural network layer, with **sparse connectivity** (many weights are set to zero) and **shared weights**

**Pooling and Sub-sampling:** replace the output with a summary statistic of the nearby outputs, e.g. max-pooling (allows invariance to small translations in the input).

53

LeCun et al, Krizhevsky et al.

## Statistical Machine Learning Hilary Term 2021

François Caron  
Department of Statistics  
University of Oxford

Slide credits and other course material can be found at:

<https://canvas.ox.ac.uk/courses/65441>

# Decision Trees

2

**Decision Trees**

Example: NHS Direct Self-help Guide (pre-covid era)

**Colds and flu**

This advice is suitable for children and adults.

Are you developing a rash that does not fade when you press a gas number or finger against it? **Dial 999**

Are you suffering from a stiff neck, headache and do your feel light hurt your eyes and/or you feel very sleepy and confused? **Dial 999**

Is there sneezing, a runny nose, a mild temperature, a sore throat, and general aches and pains? **Self-care**

It could be a common cold which antibiotics cannot treat effectively. Unless the person is very old, frail, or has other health problems, you **do not need to see your doctor**. Take paracetamol or ibuprofen (not paracetamol oral suspension, available from pharmacists), warm soups and soups. **Ask your pharmacist for advice.**

Are you feeling flushed, hot and sweaty? Do you have a high temperature (over 38°C or 100.4°F), a headache, as well as a runny nose and general aches and pains? **Self-care**

It could be this, which is generally worse than the common cold but it's helped with antibiotics. Paracetamol or ibuprofen (not paracetamol oral suspension, available from pharmacists) warm soups and soups. If your symptoms which are severe or do not go away need to see **NHS Direct**. However, if you are experiencing it is painful to bend your neck, or it tight hurts your eyes, call **NHS Direct**.

**Self-care advice**

- Take some paracetamol such as paracetamol oral suspension (ask your pharmacist), this will help to bring you down if you are having a headache, as well as a runny nose and general aches and pains.
- Increase how much fluid you or they drink.
- Some people find that a simple cough medicine helps to sooth a tickly cough.
- Flu vaccination for people who are at risk of complications. This may include the elderly, people with chronic illnesses such as heart, kidney

for example, diabetes, asthma and

# Classification and Regression Trees (CART)

Denote input domain by  $\mathcal{X}$  and let the output domain be  $\mathcal{Y} = \{1, \dots, K\}$  (classification) or  $\mathcal{Y} = \mathbb{R}$  (regression).

A decision tree gives a partition of  $\mathcal{X}$  into  $R$  disjoint sets (regions)  $\mathcal{P} = \{\mathcal{R}_1, \dots, \mathcal{R}_R\}$ , such that the fitted decision function is constant on each region  $\mathcal{R}_j \subset \mathcal{X}$ ,  $j = 1, \dots, R$ , i.e.

$$\bigcup_{j=1}^R \mathcal{R}_j = \mathcal{X} \quad h(x) = \beta_j, \quad \forall x \in \mathcal{R}_j.$$

Main strengths: easy to use, easy to interpret.

Often serve as a starting point for powerful model combination and ensemble techniques: bagging, random forests

3

**Decision Trees**

Example: NHS Direct Self-help Guide (pre-covid era)

Are you developing a rash that does not fade when you press a glass tumbler or finger against it?

yes **Emergency ("Dial 999")**  
no **Self-care**

Are you suffering from a stiff neck, headache and do you find the light hurts your eyes and/or you feeling very sleepy and confused?

yes **Emergency ("Dial 999")**  
no **Self-care**

Is there sneezing, a runny nose, a mild temperature, a sore throat, and general aches and pains?

yes **Self-care**  
no **Self-care (basic)**

Are you feeling flushed, hot and sweaty? Do you have a high temperature (over 38°C or 100.4°F), a headache, as well as a runny nose and general aches and pains?

4

5

## Decision Trees

A decision tree is a hierarchically organised structure, with each node splitting the data space into regions based on value of a single feature (attribute).

Some terminology:

Parent of a node  $c$  is the node with an arrow pointing into  $c$ .

Children of a node  $c$  are those nodes which have node  $c$  as a parent.

Root node is the top node of the tree; the only node without parents.

Leaf nodes are nodes which do not have children.

Stumps are trees with just the root node and two leaf nodes.

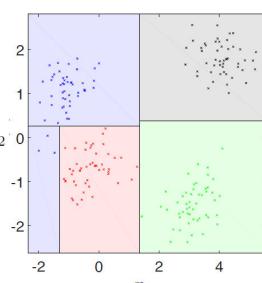
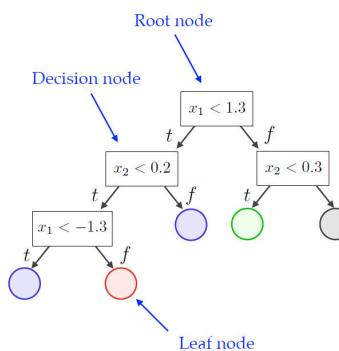
The depth of a tree is the maximal length of a path from the root node to a leaf node.

Partition of  $\mathcal{X}$  into  $R$  disjoint sets ( $\mathcal{R}_1, \dots, \mathcal{R}_R$ ) is determined by the leaves of the tree.

On each region  $\mathcal{R}_j$ , the same decision/prediction is made:  $h(x) = \beta_j$  for all  $x \in \mathcal{R}_j$  - typically as a majority vote of the data items associated to that leaf (classification under the 0-1 loss) or as their mean (regression under the squared loss)

6

## Decision tree

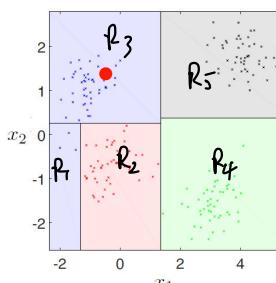
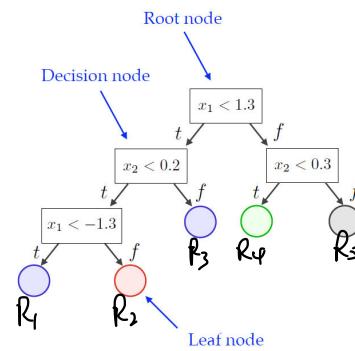


$$\mathcal{X} \in \mathbb{R}^2$$

$$\gamma = \{1, 2, 3, 4\}$$

7

## Decision tree

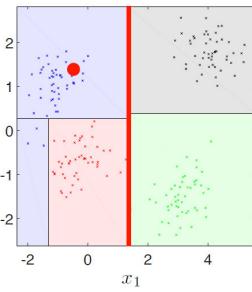
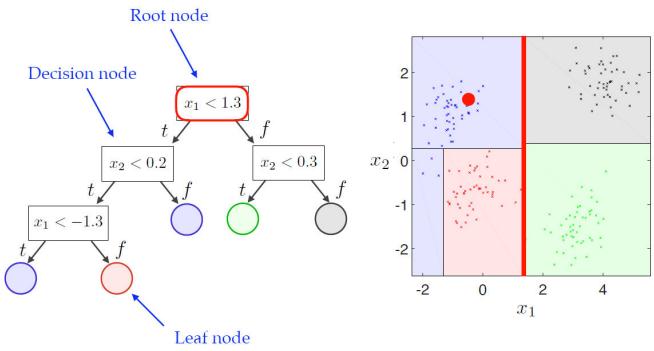


$$\bigcup_{k=1}^5 R_k = \mathcal{X}$$

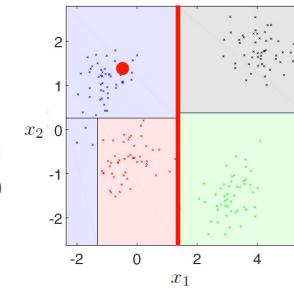
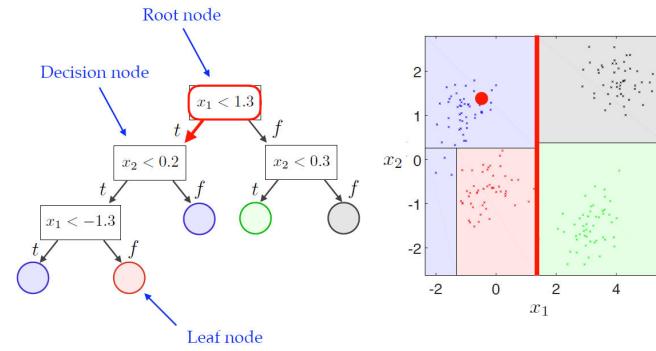
8

8

## Decision tree

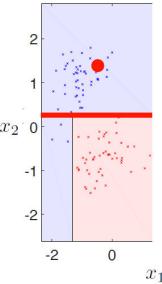
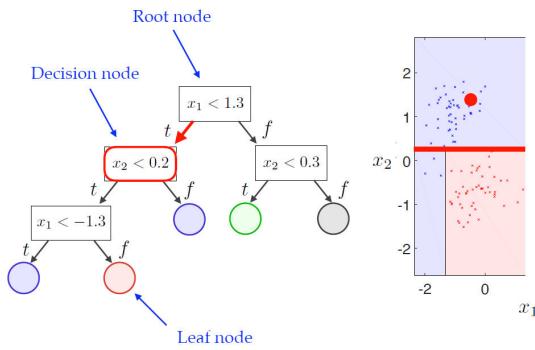


## Decision tree



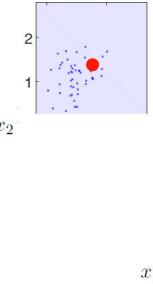
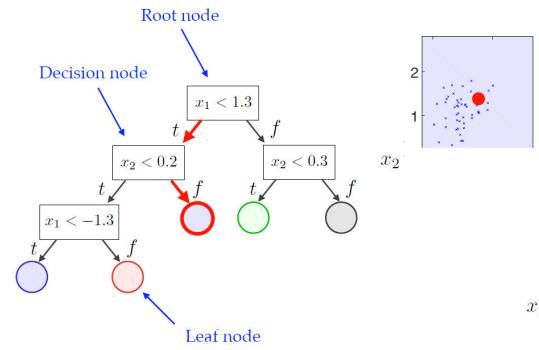
8

## Decision tree



8

## Decision tree



8

8

## Example: Iris Data

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
4.4	3.2	1.3	0.2	setosa
5.9	3.0	5.1	1.8	virginica
6.3	3.3	6.0	2.5	virginica
5.3	3.7	1.5	0.2	setosa
5.5	2.5	4.0	1.3	versicolor
6.1	2.9	4.7	1.4	versicolor
6.1	3.0	4.9	1.8	virginica
5.7	2.8	4.5	1.3	versicolor
5.4	3.0	4.5	1.5	versicolor
4.8	3.4	1.6	0.2	setosa
4.6	3.1	1.5	0.2	setosa
4.9	3.1	1.5	0.2	setosa
6.4	2.9	4.3	1.3	versicolor
.....				

Previously seen Iris data set gives the measurements in centimeters of the variables sepal length and width and petal length and width, respectively, for 50 flowers from each of 3 species of iris. The species are Iris setosa, versicolor, and virginica.

9

## Parameter estimation

The prediction rule is of the form

$$h(x) = \sum_{j=1}^R \beta_j \mathbb{1}_{x \in \mathcal{R}_j},$$

Parameters: Partition  $\{\mathcal{R}_1, \dots, \mathcal{R}_R\}$  of  $\mathcal{X}$  and weights  $\beta_1, \dots, \beta_R$   
Given the partition, estimating the parameters  $\beta_j$  via ERM is easy  
For regression, under the squared loss, we have

$$\hat{\beta}_j = \frac{\sum_{i=1}^n y_i \mathbb{1}_{x_i \in \mathcal{R}_j}}{\sum_{i=1}^n \mathbb{1}_{x_i \in \mathcal{R}_j}}$$

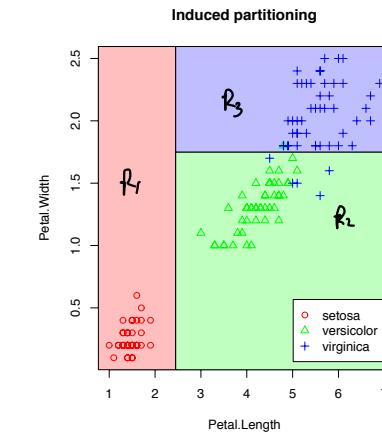
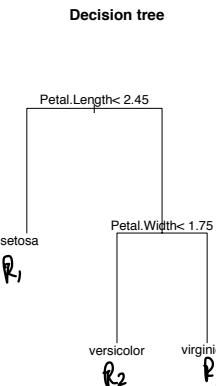
For classification problems, under the 0-1 loss, we do a majority vote within the set  $\mathcal{R}_j$

$$\hat{\beta}_j = \arg \max_{k=1, \dots, K} \sum_{i=1}^n \mathbb{1}_{y_i=k} \mathbb{1}_{x_i \in \mathcal{R}_j}$$

These estimates can be regularised as well.

11

## Example: Iris Data



Partition of  $\mathcal{X}$  into  $R$  disjoint sets  $(\mathcal{R}_1, \dots, \mathcal{R}_R)$  is determined by the leaves of the tree.

10

## Partition Estimation

Ideally, would like to find partition that achieves minimal risk: lowest mean-squared error for regression or misclassification rate for classification.

Number of potential partitions is too large to search exhaustively.

'Greedy' search heuristics for a good partition:

- Start at root.
- Determine best feature and value to split.
- Recurse on children of node.
- Stop at some point.

12

## Growth Heuristic for Regression Trees

Start with  $\mathcal{X} = \mathbb{R}^p$ .

For each feature  $j = 1, \dots, p$ , and for each value  $v \in \mathbb{R}$  that we can split on:

Split data set:

$$I_< = \{i : x_{ij} < v\}$$

$$I_> = \{i : x_{ij} \geq v\}$$

Estimate parameters:

$$\beta_< = \frac{\sum_{i \in I_<} y_i}{|I_<|}$$

$$\beta_> = \frac{\sum_{i \in I_>} y_i}{|I_>|}$$

Compute the **quality of split**, e.g., the square loss:

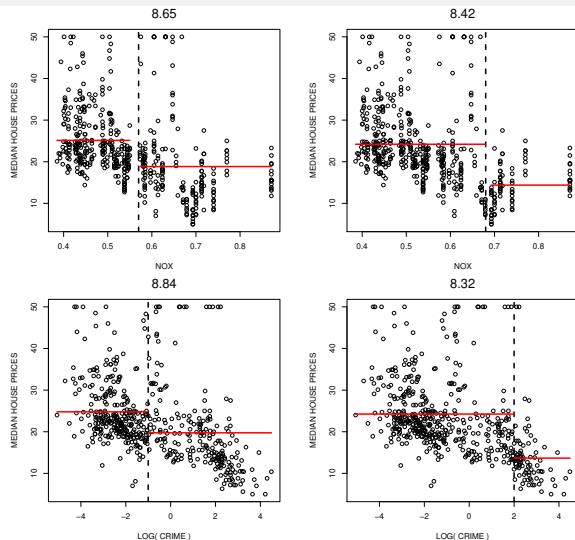
$$\sum_{i \in I_<} (y_i - \beta_<)^2 + \sum_{i \in I_>} (y_i - \beta_>)^2$$

Choose split, i.e., feature  $j$  and value  $v$ , with minimal loss.

Recurse on both children, with datasets  $(x_i, y_i)_{i \in I_<}$  and  $(x_i, y_i)_{i \in I_>}$ .

13

## Boston Housing Data



15

## Boston Housing Data

$P=13$

crim	per capita crime rate by town
zn	proportion of residential land zoned for lots over 25,000 sq.ft
indus	proportion of non-retail business acres per town
chas	Charles River dummy variable
nox	nitric oxides concentration (parts per 10 million)
rm	average number of rooms per dwelling
age	proportion of owner-occupied units built prior to 1940
dis	weighted distances to five Boston employment centres
rad	index of accessibility to radial highways
tax	full-value property-tax rate per USD 10,000
ptratio	pupil-teacher ratio by town
b	$1000(B - 0.63)^2$ where B is the proportion of blacks by town
lstat	percentage of lower status of the population
medv	median value of owner-occupied homes in USD 1000's

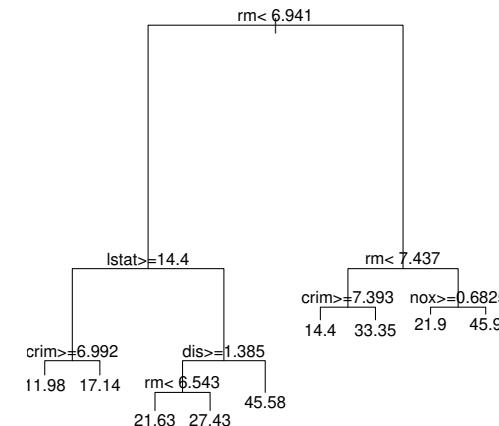
Predict median house value.

14

## Boston Housing Data

Overall, the best first split is on variable `rm`, average number of rooms per dwelling.

Final tree contains predictions in leaf nodes.



16

## Growth Heuristics for Classification Trees

For binary classification, the proportion of class 1 items in a node corresponding to a region  $\mathcal{R}$  is given by

$$\eta_1 = \frac{\sum_i \mathbb{1}_{y_i=1} \mathbb{1}_{x_i \in \mathcal{R}}}{\sum_i \mathbb{1}_{x_i \in \mathcal{R}}}$$

A split is good if both sides are more pure, i.e.  $\eta_1$  is closer to 0 or 1.

Different measures of node impurity:

- Misclassification error:  $1 - \max\{\eta_1, 1 - \eta_1\}$ .
- Gini impurity:  $2\eta_1(1 - \eta_1)$ .
- Entropy:  $-\eta_1 \log \eta_1 - (1 - \eta_1) \log(1 - \eta_1)$ .

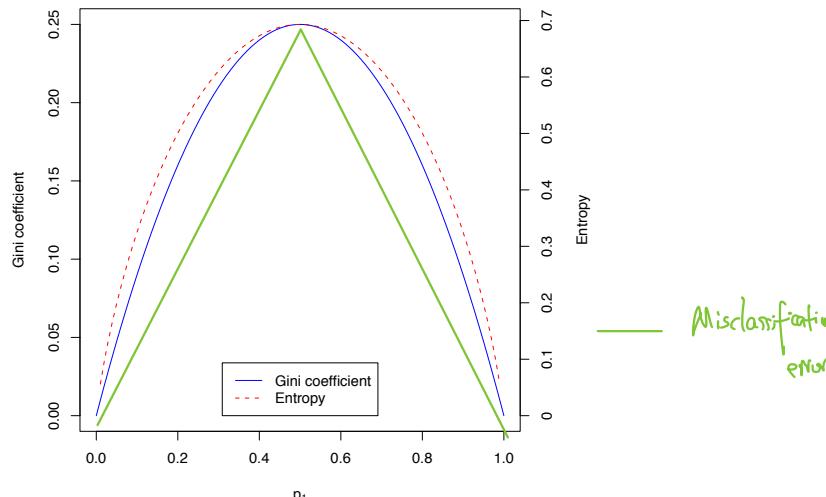
Gini and entropy preferred: differentiable and produce purer nodes.

Extension to multi-class:

- Misclassification error:  $1 - \max_k \eta_k$ .
- Gini impurity:  $\sum_{k=1}^K \eta_k(1 - \eta_k)$ .
- Entropy:  $-\sum_{k=1}^K \eta_k \log \eta_k$ .

Stops once a node has insufficient number of items, or is pure.

17



Misclassification error?

## Growth Heuristics for Classification Trees

Consider splitting  $\mathcal{R}$  into two non-overlapping regions  $\mathcal{R}_l$  and  $\mathcal{R}_r$ , where  $\mathcal{R}_l \cup \mathcal{R}_r = \mathcal{R}$

Denote  $q_l = \frac{\#\{i | x_i \in \mathcal{R}_l\}}{\#\{i | x_i \in \mathcal{R}\}}$  the proportion of the examples in  $\mathcal{R}$  assigned to the left partition and  $q_r = 1 - q_l$  the proportion of the examples in  $\mathcal{R}$  assigned to the right partition

Denote  $\eta_1, \eta_{1,l}, \eta_{1,r}$  the proportion of class 1 items respectively in regions  $\mathcal{R}, \mathcal{R}_l$  and  $\mathcal{R}_r$

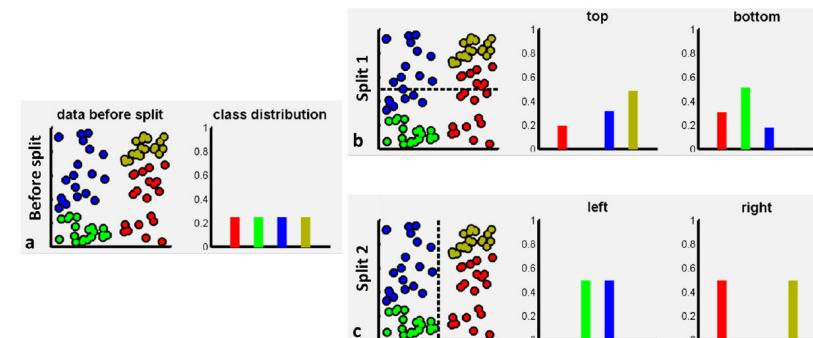
Let  $i(\eta_1), i(\eta_{1,l}), i(\eta_{1,r})$  the corresponding impurity measures

We choose the split that maximises the change in the impurity measure

$$i(\eta_1) - q_l i(\eta_{1,l}) - q_r i(\eta_{1,r})$$

18

## Growth Heuristics for Classification Trees



Misclassification error is 0.25 for both splits.

20

Adapted from Criminisi et al., Decision Forests, 2012.

## Example: Pima Indians Diabetes Dataset

The subjects: women who were at least 21 years old, of Pima Indian heritage living near Phoenix, Arizona.

Tested for diabetes according to World Health Organisation criteria.

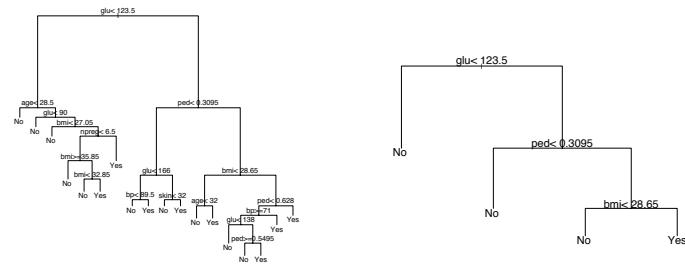
## Features:

- number of pregnancies (npreg)  
plasma glucose concentration (glu)  
diastolic blood pressure (bp)  
tricep skin fold thickness (skin)  
body mass index(bbi)  
diabetes pedigree function (ped)  
age (age)

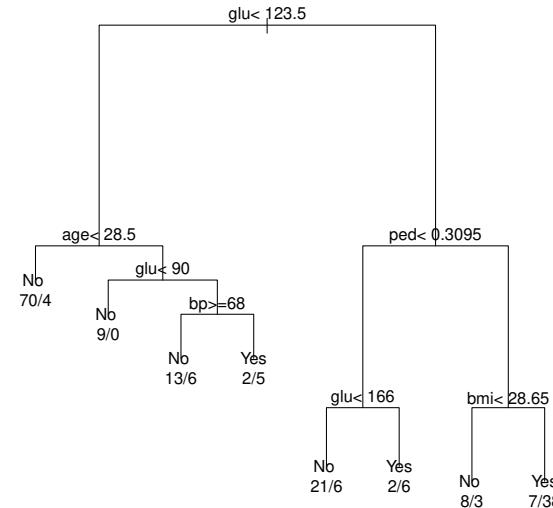
21

## Example: Pima Indians Diabetes Dataset

## Two possible trees.



## Example: Pima Indians Diabetes Dataset



22

## Model Complexity

When should tree growing be stopped?

Will need to control complexity to prevent overfitting, and in general find optimal tree size with best predictive performance.

## A regularised objective

$$\widehat{R}_n(h_T) + \frac{\lambda}{n} \times \text{size}(T)$$

with  $\text{size}(T)$  the number of leaf nodes.

Grow the tree from scratch and stop once the criterion objective starts to increase.

First grow the full tree and prune nodes (starting at leaves), until the objective starts to increase. ~~cut back~~

Second option is preferred as the choice of tree is less sensitive to “wrong” choices of split points and variables to split on in the first stages of tree fitting.

Use a validation set or cross-validation to determine optimal  $\lambda$ .

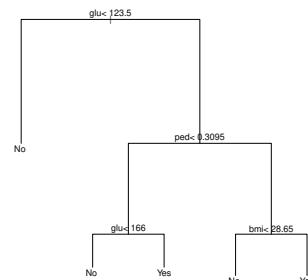
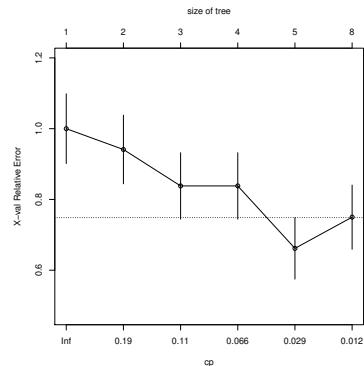
23

24

## Model Complexity

Pima Indian dataset

Cross-validation risk versus regularisation parameter  $\lambda$



25

## Summary

### Advantages

- Interpretable predictions
- Handles both numerical and categorical data

### Disadvantages

- Instability: a small change in the data may result in a very different series of splits
- Lack of smoothness for regression
- Often poor predictive performances

## Computational Considerations

### Numerical input features $x_{ij}$

We could split on any feature, with any threshold

However, for a given feature, the only split points we need to consider are the  $n$  values in the training data for this feature.

If we sort each feature by these  $n$  values, we can quickly compute our metric of interest (squared loss or impurity)

This takes  $O(pn \log n)$  time

### Categorical input features $x_{ij}$

Assuming  $q$  distinct categories, there are  $2^{q-1} - 1$  possible partitions we can consider.

26

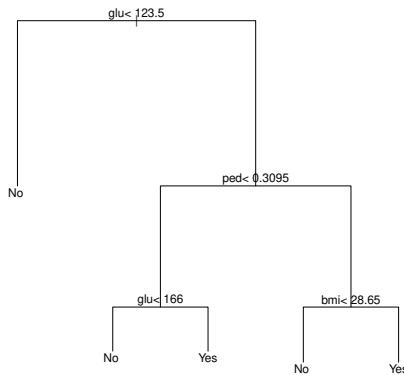
## Bagging

27

28

## Model Variability

Bagging



Is the tree 'stable' if training data were slightly different?

29

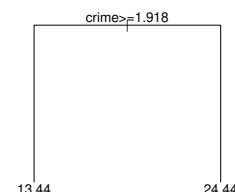
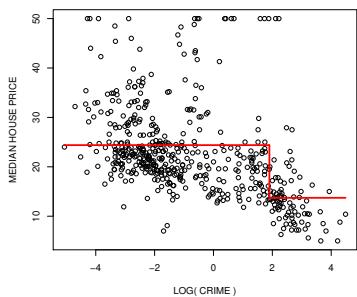
## Bootstrap for Regression Trees

Bagging

Regression for Boston housing data.

Predict median house prices based only on crime rate.

Use decision stump—the simplest tree with a single split at root.



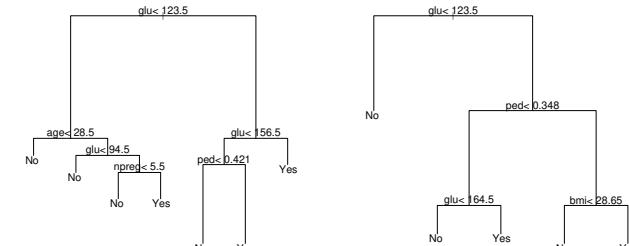
31

## Bootstrap for Classification Trees

Bagging

The bootstrap is a way to assess the variance of estimators.

Fit multiple trees, each on a bootstrapped sample. This is a data set obtained by sampling with replacement  $n$  times from training set.



$$\left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \right) \xrightarrow{\text{Bootstrap}} \left( \begin{array}{c} x'_1 \\ x'_2 \\ x'_3 \\ x'_4 \end{array} \right), \left( \begin{array}{c} x''_1 \\ x''_2 \\ x''_3 \\ x''_4 \end{array} \right), \dots$$

30

## Bootstrap for Regression Trees

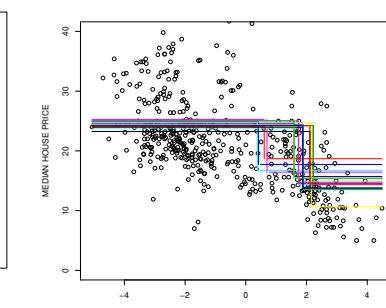
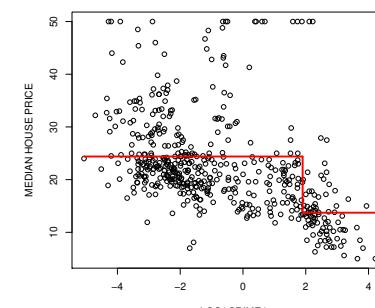
Bagging

We fit a predictor  $\hat{h}^{(d)}(x)$  on the data  $d = \{(x_i, y_i)\}_{i=1}^n$ .

Assess the variance of  $\hat{h}^{(D)}(x)$  by taking  $B = 20$  bootstrap samples of the original data, and obtaining bootstrap estimators

$$\hat{h}^b(x), \quad b = 1, \dots, B$$

Each tree  $\hat{h}^b$  is fitted on the resampled data  $(x_{j_i}, y_{j_i})_{i=1}^n$  where each  $j_i$  is chosen randomly from  $\{1, \dots, n\}$  with replacement.



32

## Bagging

Bagging

Bagging (Bootstrap Aggregation): average across all  $B$  trees fitted on different bootstrap samples.

For  $b = 1, \dots, B$ :

Draw indices  $(j_1, \dots, j_n)$  from the set  $\{1, \dots, n\}$  with replacement.  
Fit the model, and form predictor  $\hat{h}^b(x)$  based on bootstrap sample

$$(x_{j_1}, y_{j_1}), \dots, (x_{j_n}, y_{j_n})$$

Form bagged estimator

$$\hat{h}_{Bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{h}^b(x)$$

33

## Variance Reduction in Bagging

Bagging

Suppose, in an ideal world, our estimators  $\hat{h}^b$  are each based on different independent datasets of size  $n$  from the true joint distribution of  $X, Y$ .

The aggregated estimator would then be

$$\hat{h}_{ag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{h}^b(x) \rightarrow \bar{h}(x) = \mathbb{E}_D[\hat{h}(x)] \quad \text{a.s. as } B \rightarrow \infty$$

LLN

where expectation is with respect to datasets of size  $n$ .

The conditional risk under the squared-loss is:

$$\begin{aligned} & \mathbb{E}_D[(Y - \hat{h}_{ag}(X))^2 | X = x] \xrightarrow{\text{Bias}} \\ & = \mathbb{E}_D[(Y - \bar{h}(X))^2 | X = x] + \mathbb{E}_D[(\bar{h}(X) - \hat{h}_{ag}(X))^2 | X = x] \xrightarrow{\text{Variance}} \\ & \rightarrow \mathbb{E}_D[(Y - \bar{h}(X))^2 | X = x] \quad \text{as } B \rightarrow \infty. \end{aligned}$$

Aggregation reduces the squared loss by eliminating variance of  $\hat{h}(x)$ .

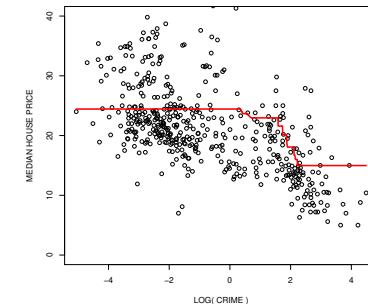
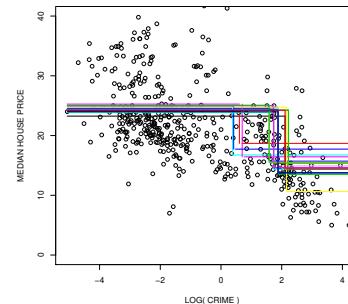
In bagging, the bootstrap samples are not independent, but variance reduction still applies at the cost of a small increase in bias.

Bagging is most useful for flexible estimators with high variance (and low bias).

35

## Bagging

Bagging



Bagging smooths out the drop in the estimate of median house prices.  
Bagging reduces the variance of predictions, i.e. when taking expectations over a random dataset  $D$ :

$$\mathbb{E}_D[(\hat{h}(x) - \mathbb{E}_D[\hat{h}(x)])^2] \geq \mathbb{E}_D[(\hat{h}_{Bag}(x) - \mathbb{E}_D[\hat{h}_{Bag}(x)])^2]$$

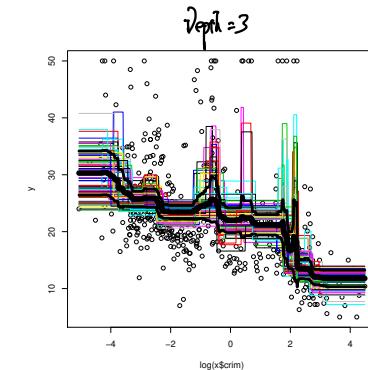
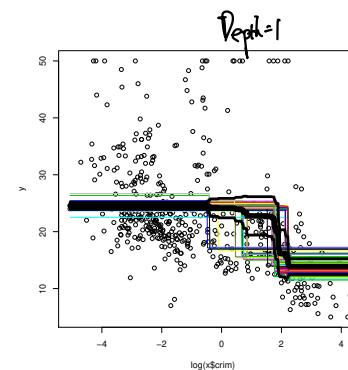
34

## Variance Reduction in Bagging

Bagging

Deeper trees have higher complexity and variance.

Compare bagging trees of depth 1 and 3.

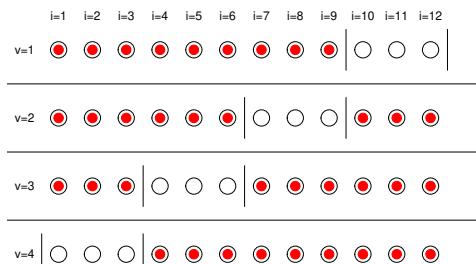


36

## Out-of-bag Test Error Estimation

How well does bagging do? Can we estimate generalisation performance, and tune hyperparameters?

Answer 1: cross-validation.



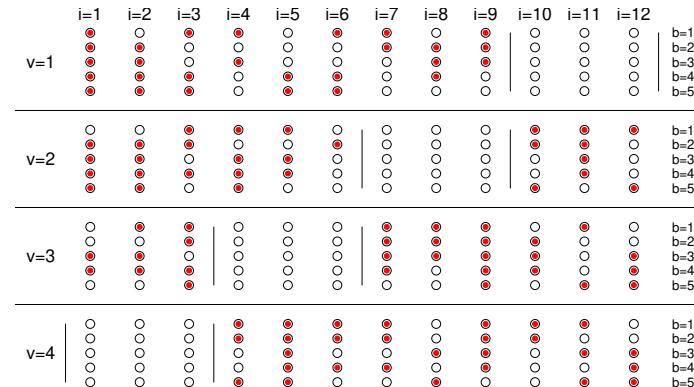
For each  $v = 1, \dots, V$ ,  
fit  $\hat{h}_{Bag}$  on the training samples.  
predict on validation set.

Compute the CV error by averaging the loss across all test observations.

37

## Out-of-bag Test Error Estimation

But to fit  $\hat{h}_{Bag}$  on the training set for each  $v = 1, \dots, V$ , we have to train on  $B$  bootstrap samples!

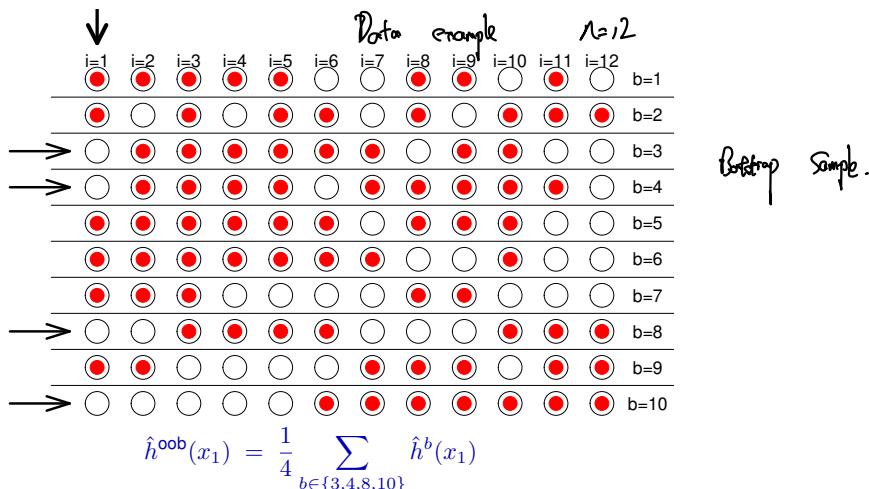


Answer 2: Out-of-bag test error estimation.

38

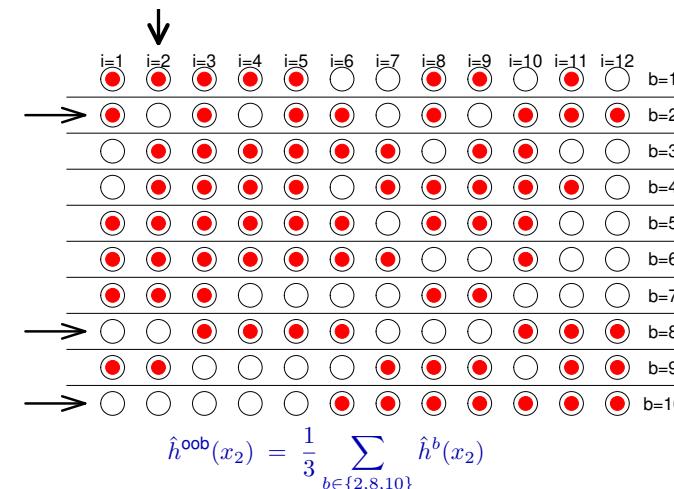
## Out-of-bag Test Error Estimation

Idea: test on the “unused” data points in each bootstrap iteration to estimate the test error.



39

40



## Out-of-bag Test Error Estimation

For each  $i = 1, \dots, n$ , the out-of-bag sample is:

$$\tilde{B}_i = \{b : x_i \text{ is not in training set}\} \subseteq \{1, \dots, B\}.$$

Construct the out-of-bag estimate at  $x_i$ :

$$\hat{h}^{\text{oob}}(x_i) = \frac{1}{|\tilde{B}_i|} \sum_{b \in \tilde{B}_i} \hat{h}^b(x_i)$$

Out-of-bag risk:

$$\hat{R}^{\text{oob}} = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{h}^{\text{oob}}(x_i))$$

41

## Example: Boston Housing Dataset

Apply out of bag test error estimation to select optimal tree depth and assess performance of bagged trees for Boston Housing data.

Use the entire dataset with  $p = 13$  predictor variables.

## Out-of-bag Test Error Estimation

We need  $|\tilde{B}_i|$  to be reasonably large for all  $i = 1, \dots, n$ .

The probability  $\pi^{\text{oob}}$  of an observation NOT being included in a bootstrap sample  $(j_1, \dots, j_n)$  (and hence being 'out-of-bag') is:

$$\pi^{\text{oob}} = \prod_{i=1}^n \left(1 - \frac{1}{n}\right) \xrightarrow{n \rightarrow \infty} \frac{1}{e} \approx 0.367.$$

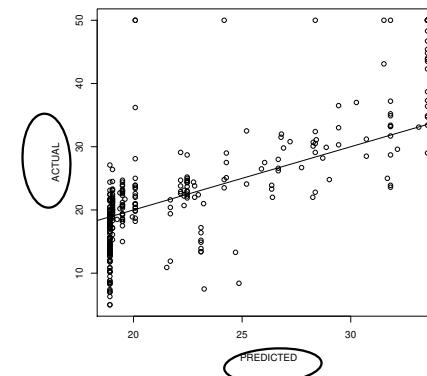
Hence  $\mathbb{E}[|\tilde{B}_i|] \approx 0.367B$

In practice, number of bootstrap samples  $B$  is typically between 200 and 1000, meaning that the number  $|\tilde{B}_i|$  of out-of-bag samples will be approximately in the range 70 – 350.

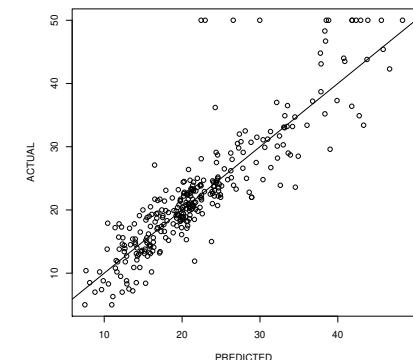
42

## Example: Boston Housing Dataset

For depth  $d = 1$ .



For depth  $d = 10$ .



43

44

## Example: Boston Housing Dataset

Test error as a function of tree depth  $d$ :

tree depth $d$	1	2	3	4	5	10	30
single tree $\hat{h}$	60.7	44.8	32.8	31.2	27.7	26.5	27.3
bagged trees $\hat{h}_{Bag}$	43.4	27.0	22.8	21.5	20.7	20.1	20.1

Without bagging, the optimal tree depth seems to be  $d = 10$ .

With bagging, we could also take the depth up to  $\underline{d = 30}$ .

45

## Random Forests

## Summary

Bagging reduces variance and prevents overfitting

Often improves accuracy in practice.

Bagged trees cannot be displayed as nicely as single trees and some of the interpretability of trees is lost.

46

## Random Forests and Extremely Randomized Trees

 Random forests are similar to bagged decision trees with a few key differences:

For each split point, the search is not over all  $p$  variables but just over some  $p_{max} \leq p$  randomly chosen ones (where e.g.  $p_{max} = \lfloor \sqrt{p} \rfloor$ )

No pruning necessary. Trees can be grown until each node contains just very few observations (1 or 5).

Random forests tend to produce better predictions than bagging.

Results often not sensitive to the tuning parameter  $p_{max}$ .

Even more random methods, e.g. extremely randomized trees:

For each split point, sample  $p_{max}$  variables each with a random value to split on, and pick the best one.

Often works even when  $p_{max}$  equals 1!

Often produce very accurate predictions, and amongst the top performing methods in machine learning competitions.

47

48

## Random Forests

TABLE 2  
Test set misclassification error (%)

Data set	Forest	Single tree
Breast cancer	2.9	5.9
Ionosphere	5.5	11.2
Diabetes	24.2	25.3
Glass	22.0	30.4
Soybean	5.7	8.6
Letters	3.4	12.4
Satellite	8.6	14.8
Shuttle $\times 10^3$	7.0	62.0
DNA	3.9	6.2
Digit	6.2	17.1

From Breiman, Statistical Modelling: the two cultures, 2001.

49

## Variable “Importance”

Tree ensembles have better performance, but decision trees are more interpretable.

How to interpret a forest of trees ?

Idea: denote by  $\hat{e}$  the out-of bag estimate of the loss when using the original data samples. For each variable  $j \in \{1, \dots, p\}$ ,

permute randomly the  $j$ -th predictor variable to generate a new set of samples  $(\tilde{X}_1, Y_1), \dots, (\tilde{X}_n, Y_n)$ , i.e.,  $\tilde{X}_{ij} = X_{\tau(i)j}$ , for a permutation  $\tau$ .

compute the out-of-bag estimate  $\hat{e}_j$  of the prediction error with these new samples.

A measure of importance of variable  $j$  is then  $\hat{e}_j - \hat{e}$ , the increase in error rate due to a random permutation of the  $j$ -th variable.

## Random Forests

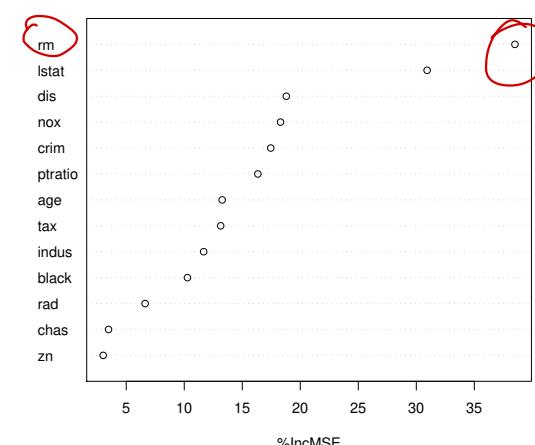
Comparison of 179 classifiers on 121 datasets. Random forests come top with SVMs close behind.

Rank	Acc.	$\kappa$	Classifier
<b>32.9</b>	82.0	63.5	parRF_t (RF)
33.1	<b>82.3</b>	<b>63.6</b>	rf_t (RF)
36.8	81.8	62.2	svm_C (SVM)
38.0	81.2	60.1	svmPoly_t (SVM)
39.4	81.9	62.5	rforest_R (RF)
39.6	82.0	62.0	elm_kernel_m (NNET)
40.3	81.4	61.1	svmRadialCost_t (SVM)
42.5	81.0	60.0	svmRadial_t (SVM)
42.9	80.6	61.0	C5.0_t (BST)
44.1	79.4	60.5	avNNet_t (NNET)

From Delgado et al, 2014

50

Example for Boston Housing data.



51

52

## Random Forests

### Ensemble Methods

Bagging and random forests are examples of **ensemble methods**, where predictions are based on an ensemble of many individual predictors.

Many other ensemble learning methods: boosting, stacking, mixture of experts, Bayesian model averaging etc.

Often gives significant boost to predictive performance.

53



## Random Forests

### Summary

#### Advantages of random forests

- Easy to use
- Fast
- State-of-the-art for many application
- Particularly good for small/medium size datasets
- Requires little tuning

#### Disadvantages

- Typically worse than deep learning on huge datasets
- Limited Interpretability

54



## Statistical Machine Learning

### Hilary Term 2021

François Caron  
Department of Statistics  
University of Oxford

Slide credits and other course material can be found at:

<https://canvas.ox.ac.uk/courses/65441>

## Boosting

1

2

## Boosting

Binary classification

$$\gamma = \mathbb{F}_{1,1}$$

Weak classifiers  $h_t : \mathcal{X} \rightarrow \{-1, 1\}$  in some base hypothesis class  $\mathcal{H}$ , for  $t = 1, \dots, T$

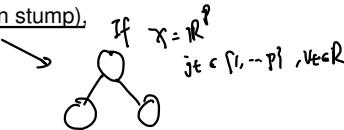
Simple classifiers with few parameters to learn (e.g. decision stump), typically with high bias and small variance

Boosting prediction rule

$$h(x) = \text{sign}(f(x))$$

where

$$f(x) = \sum_{t=1}^T \beta_t h_t(x).$$



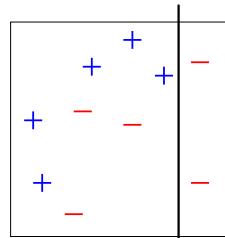
Weighted combination of weak classifiers

Parameters to estimate:  $(\beta_t, h_t)_{t=1, \dots, T}$

3

## Example

$n = 10$  data examples and  $p = 2$  features



The data points are clearly not linearly separable

In the beginning, all data points have equal weights (weights are indicated by the size of the data markers "+" or "-")

Base class  $\mathcal{H}$  is the class of decision stumps: trees with a single split and two leaf nodes, classifying data based on a single attribute ( $x_{i1}$  or  $x_{i2}$ )

5



## Adaboost Algorithm

Dataset  $d = \{(x_i, y_i)\}$ , where  $y_i \in \{+1, -1\}$

Initialise weights  $\bar{w}_{i,1} = \frac{1}{n}$  for every training sample  $i = 1, \dots, n$

For  $t = 1$  to  $T$

Train the weak classifier  $h_t$  using current weights  $\bar{w}_{i,t}$ , by minimising

$$\hat{h}_t = \arg \min_{h_t \in \mathcal{H}} R_{\bar{w}_t}(h_t)$$

where  $R_{\bar{w}_t}(h_t) = \sum_{i=1}^n \bar{w}_{i,t} \mathbb{1}_{y_i \neq h_t(x_i)}$  is the weighted classification error.

Compute contribution for this classifier:  $\hat{\beta}_t = \frac{1}{2} \ln \frac{1 - \hat{\epsilon}_t}{\hat{\epsilon}_t}$  where  $\hat{\epsilon}_t = R_{\bar{w}_t}(\hat{h}_t)$ .

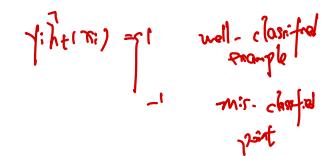
Update weights on training points

$$\bar{w}_{i,t+1} \propto \bar{w}_{i,t} e^{-\hat{\beta}_t y_i \hat{h}_t(x_i)}$$

and normalise them such that  $\sum_i \bar{w}_{i,t+1} = 1$ .

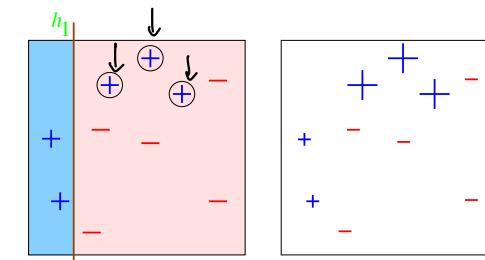
Output the final classifier

$$\hat{h}(x) = \text{sign} \left( \sum_{t=1}^T \hat{\beta}_t \hat{h}_t(x) \right)$$



4

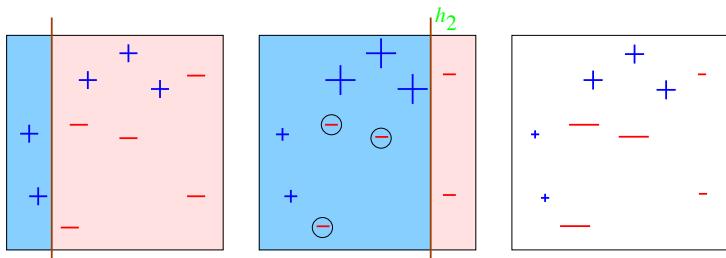
## Round 1: $t = 1$



3 misclassified (with circles):  $\hat{\epsilon}_1 = 0.3 \rightarrow \hat{\beta}_1 = 0.42$ .

Weights recomputed; the 3 misclassified data points receive larger weights

6

Round 2:  $t = 2$ 

3 misclassified (with circles):  $\hat{\epsilon}_2 = 0.21 \rightarrow \hat{\beta}_2 = 0.65$ .

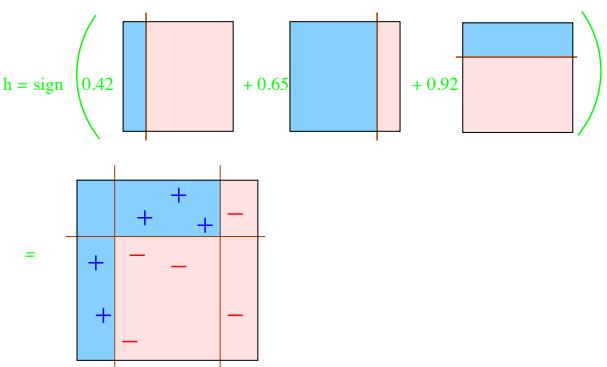
Note that  $\hat{\epsilon}_2 \neq 0.3$  as those 3 data points have weights less than 1/10

3 misclassified data points get larger weights

Data points classified correctly in both rounds have small weights

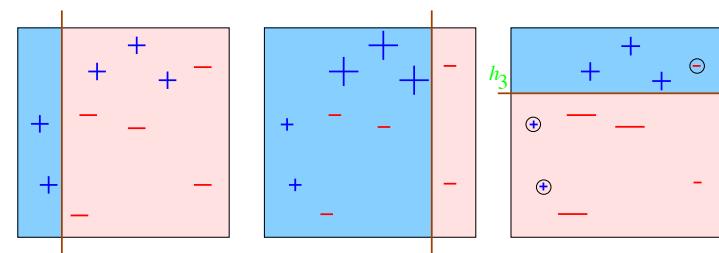
7

## Final classifier: combining 3 classifiers



All data points are now classified correctly!

9

Round 3:  $t = 3$ 

3 misclassified (with circles):  $\hat{\epsilon}_3 = 0.14 \rightarrow \hat{\beta}_3 = 0.92$ .

Previously correctly classified data points are now misclassified, hence our error is low; what's the intuition?

Since they have been consistently classified correctly, this round's mistake will hopefully not have a huge impact on the overall prediction

8

## Why AdaBoost works?

Classifier

$$h(x) = \text{sign}(f(x)) = \begin{cases} 1 & \text{if } f(x) > 0 \\ -1 & \text{if } f(x) < 0 \end{cases}$$

where  $f$  is the discriminant function.

Under the 0-1 loss

$$L(y, h(x)) = \begin{cases} 0 & \text{if } y f(x) > 0 \\ 1 & \text{if } y f(x) < 0 \end{cases}$$

The discriminant function  $f(x)$  and the target label  $y$  should have the same sign to avoid a loss of 1.

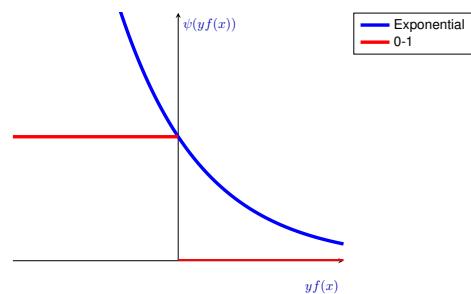
10

## Surrogate Exponential loss

As we discussed in the lectures on linear models, the  $0 - 1$  loss is non-convex and difficult to optimise.

### Surrogate Exponential Loss

$$\tilde{L}(y, h(x)) = \psi(yf(x)) = e^{-yf(x)}$$



11

We have

$$\sum_{i=1}^n \exp(-y_i (\hat{f}_{t-1}(x_i) + \beta_t h_t(x_i))) = \sum_{i=1}^n w_{i,t} \exp(-y_i \beta_t h_t(x_i))$$

where  $w_{i,t} = e^{-y_i \hat{f}_{t-1}(x_i)}$ .

Note that the Adaboost weights are  $\bar{w}_{i,t} = w_{i,t}/(\sum_j w_{j,t})$ .

We have  $y_i h_t(x_i) \in \{1, -1\}$ , hence

$$\begin{aligned} \sum_{i=1}^n w_{i,t} \exp(-y_i \beta_t h_t(x_i)) &= \sum_{i=1}^n w_{i,t} e^{\beta_t} \mathbb{1}_{y_i \neq h_t(x_i)} + \sum_{i=1}^n w_{i,t} e^{-\beta_t} \mathbb{1}_{y_i = h_t(x_i)} \\ &= (e^{\beta_t} - e^{-\beta_t}) \sum_{i=1}^n w_{i,t} \mathbb{1}_{y_i \neq h_t(x_i)} + e^{-\beta_t} \sum_{i=1}^n w_{i,t} \\ &= \left( \sum_{i=1}^n w_{i,t} \right) \left( (e^{\beta_t} - e^{-\beta_t}) \sum_{i=1}^n \bar{w}_{i,t} \mathbb{1}_{y_i \neq h_t(x_i)} + e^{-\beta_t} \right) \end{aligned}$$

Minimising wrt  $h_t$  gives Step 1 of Adaboost

$$\hat{h}_t = \arg \min_{h_t \in \mathcal{H}} \sum_{i=1}^n \bar{w}_{i,t} \mathbb{1}_{y_i \neq h_t(x_i)}$$

13

## ERM under Surrogate Exponential loss

Adaboost aims at minimising the empirical risk under the surrogate exponential loss

$$\frac{1}{n} \sum_{i=1}^n e^{-y_i f(x_i)} = \frac{1}{n} \sum_{i=1}^n e^{-y_i (\sum_{t=1}^T \beta_t h_t(x_i))}$$

Direct minimisation is not possible, and Adaboost uses a stagewise/greedy approach: (forward stagewise additive modelling)

Initialise  $\hat{f}_0(x) = 0$ .

At iteration  $t = 1, \dots, T$ , add a new weighted weak learner into the model:

$$\begin{aligned} (\hat{\beta}_t, \hat{h}_t) &= \arg \min_{\beta_t \in \mathbb{R}, h_t \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \exp \left( -y_i \left( \hat{f}_{t-1}(x_i) + \beta_t h_t(x_i) \right) \right) \\ \hat{f}_t(x) &= \hat{f}_{t-1}(x) + \hat{\beta}_t \hat{h}_t(x) \end{aligned}$$

Let's prove this is equivalent to the Adaboost algorithm update

12

For  $\beta_t$ , we need to minimise

$$(e^{\beta_t} - e^{-\beta_t}) \hat{e}_t + e^{-\beta_t}$$

where  $\hat{e}_t = \sum_{i=1}^n \bar{w}_{i,t} \mathbb{1}_{y_i \neq \hat{h}_t(x_i)}$ .

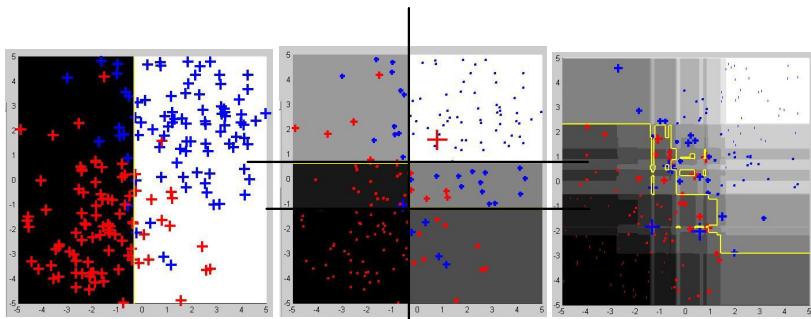
Taking the derivative and setting to 0 gives

$$\hat{\beta}_t = \frac{1}{2} \log \frac{1 - \hat{e}_t}{\hat{e}_t}$$

which is exactly Step 2 of Adaboost.

14

## AdaBoost with decision stumps



The degree of blackness represents the confidence in the red class. The size of datapoints represents their weight. Decision boundary in yellow.  
Left: After 1 iteration, Middle: After 3 iterations, Right: After 120 iterations.

15

Example from Murphy, p.560; generating script written by R.Stapenhurst

## Base learners

We can use other base learners than decision stumps, as long as we are minimising the weighted classification error

Typically use a simple, high bias classifier, whose parameters are easy to estimate



## Example: Netflix

### The Netflix Prize

<http://www.netflixprize.com>

Training data is a set of users and past ratings (1 to 5 stars).

Construct a classifier that predicts user rating for unrated movies.

Winning team (BellKor's Pragmatic Chaos) employed boosting. They received 1M\$ in 2009 for a 10.06% improvement.

16

## Other surrogate loss functions

Possible to use other surrogate loss functions  $\psi$ , aiming to minimise

$$\frac{1}{n} \sum_{i=1}^n \psi \left( y_i \left( \sum_{t=1}^T \beta_t h_t(x_i) \right) \right)$$

The  $\psi$ -boosting algorithm proceeds as follows

Initialise  $\hat{f}_0(x) = 0$ .

At iteration  $t = 1, \dots, T$ , add a new weighted weak learner into the model:

$$\begin{aligned} (\hat{\beta}_t, \hat{h}_t) &= \arg \min_{\beta_t \in \mathbb{R}, h_t \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \psi \left( y_i \left( \hat{f}_{t-1}(x_i) + \beta_t h_t(x_i) \right) \right) \\ \hat{f}_t(x) &= \hat{f}_{t-1}(x) + \hat{\beta}_t \hat{h}_t(x) \end{aligned}$$

If minimisation not tractable, perform a gradient descent step instead

If  $\psi(x) = \log(1 + e^{-x})$ , this is known as logit-boost.

17

18

## $L_2$ -Boosting

Can also use boosting for regression

Under the squared-loss, known as  $L_2$ -Boosting

For example, use regression trees as weak learner

New weak learners are obtained by fitting the residuals:

$$L \left( y_i, \hat{f}_{t-1}(x_i) + \beta_t h_t(x_i) \right) = \underbrace{(y_i - \hat{f}_{t-1}(x_i) - \beta_t h_t(x_i))^2}_{r_{i,t}} = L(y_i - \underbrace{\hat{f}_{t-1}(x_i)}_{r_{i,t}}, \beta_t h_t(x_i))$$

$L_2$ -Boosting algorithm

Initialise  $\hat{f}_0(x) = 0$ .

For  $t = 1, \dots, T$ , compute current residuals

$$\underline{r_{i,t} = y_i - \hat{f}_{t-1}(x_i)},$$

and fit the residuals  $\{(x_i, r_{i,t})\}_{i=1}^n$  to obtain the term  $\underline{\beta_t \hat{h}_t(x)}$  to be added to the expansion.

## Boosting: Summary

Boosting is a bias-reduction technique, as opposed to bagging.

Resistant to overfitting (the testing error typically stays flat for a large number of iterations - but will eventually go up).

Can be understood as functional gradient descent, leading to a generic framework called gradient boosting.

Further reading: Hastie et al, Chapter 10; Murphy, Section 16.4.