

Multivariate Statistics and Machine Learning

MATH38161

Korbinian Strimmer

9 September 2020

Contents

Preface	5
About these notes	5
About the author	5
About the module	5
Acknowledgements	6
1 Multivariate random variables	7
1.1 Why multivariate statistics?	7
1.2 Basics	8
1.3 Multivariate normal distribution	11
1.4 Estimation in large sample and small sample settings	16
1.5 Discrete multivariate distributions	23
1.6 Continuous multivariate distributions	25
2 Transformations and dimension reduction	31
2.1 Linear Transformations	31
2.2 Nonlinear transformations	33
2.3 Whitening transformations	35
2.4 Natural whitening procedures	38
2.5 Principal Component Analysis (PCA)	42
2.6 PCA correlation loadings and plot	43
2.7 CCA whitening (Canonical Correlation Analysis)	45
3 Clustering / unsupervised Learning	47
3.1 Overview of clustering	47
3.2 Hierarchical clustering	50
3.3 K-means clustering	54
3.4 Mixture models	59
3.5 Fitting mixture models to data	63
4 Classification / supervised learning	69
4.1 Introduction	69
4.2 Bayesian discriminant rule or Bayes classifier	71
4.3 Normal Bayes classifier	71
4.4 The training step — learning QDA, LDA and DDA classifiers from data	74
4.5 Goodness of fit and variable selection	76
4.6 Estimating prediction error	79

5	Multivariate dependencies	81
5.1	Measuring the association between two sets of random variables	81
5.2	Graphical models	82
6	Nonlinear and nonparametric models	89
6.1	Limits of linear models and correlation	89
6.2	Mutual information as generalised correlation	91
6.3	Nonlinear regression models	95
6.4	Random forests	96
6.5	Gaussian processes	97
6.6	Neural networks	99
A	Brief refresher on matrices	101
A.1	Matrix notation	101
A.2	Simple special matrices	102
A.3	Simple matrix operations	102
A.4	Orthogonal matrices	102
A.5	Eigenvalues and eigenvalue decomposition	103
A.6	Singular value decomposition	103
A.7	Positive (semi-)definiteness, rank, condition	103
A.8	Trace and determinant of a matrix and eigenvalues	104
A.9	Functions of matrices	104
A.10	Matrix calculus	105
B	Further study	107
B.1	Recommended reading	107
B.2	Advanced reading	107

Preface

About these notes

This is the course text for MATH38161, an introductory course in **Multivariate Statistics and Machine Learning** for third year mathematics students.

These notes will be updated from time to time. To view the current version in your browser visit the [online MATH38161 lecture notes](#). You may also [download the MATH38161 lecture notes as PDF](#).

About the author

My name is Korbinian Strimmer and I am a Professor in Statistics in the [Statistics group of the Department of Mathematics at the University of Manchester](#). You can find more information about me on [my home page](#).

I have first taught this module in the Winter term 2018 at the University of Manchester.

I hope you enjoy the course. If you have any questions, comments, or corrections then please email me at korbinian.strimmer@manchester.ac.uk

About the module

Topics covered

The MATH38161 module is designed to run over the course of 11 weeks. It has six parts, each covering a particular aspect of multivariate statistics and machine learning:

1. Multivariate random variables and estimation in large and small sample settings (W1 and W2)
2. Transformations and dimension reduction (W3 and W4)
3. Unsupervised learning/clustering (W5 and W6)
4. Supervised learning/classification (W7 and W8)
5. Measuring and modelling multivariate dependencies (W9)
6. Nonlinear and nonparametric models (W10, W11)

This module focuses on:

- *Concepts and methods* (not on theory)
- *Implementation and application in R*
- *Practical data analysis and interpretation* (incl. report writing)
- *Modern tools in data science and statistics* (R markdown, R studio)

Additional support material

For organisational details (including details on assessment) please visit the [course home page](#) (academic year 2020/21).

If you are a student and enrolled for this module you can find further information and materials on the [University of Manchester Blackboard](#).

On Blackboard and the course homepage you find these lecture notes plus:

- A study guide with a week by week plan what to read and study.
- 11 worksheets, one for each week, with examples (theory and application in R) and solutions.
- R Markdown code for the worksheets so that you can run them directly in R Studio.
- Links to the online MATH30161 lectures (visualiser style).

On Blackboard you also find the exam questions of previous years (without solution) as well as the coursework instructions.

Furthermore, there is also an [MATH38161 online reading list](#) provided at the UoM library.

Acknowledgements

Many thanks to [Beatriz Costa Gomes](#) for her help to compile the first draft of these course notes in the winter term 2018 while she was a graduate teaching assistant for this course. I also thank the many students who suggested corrections.

Chapter 1

Multivariate random variables

1.1 Why multivariate statistics?

Science uses experiments to verify hypotheses about the world. Statistics provides tools to quantify this procedure and offers methods to link data (experiments) with probabilistic models (hypotheses). Since the world is complex we need complex models and complex data, hence the need for multivariate statistics and machine learning.

Specifically, multivariate statistics (as opposed to univariate statistics) is concerned with methods and models for **random vectors** and **random matrices**, rather than just random univariate (scalar) variables. Therefore, in multivariate statistics we will frequently make use of matrix notation.

Closely related to multivariate statistics (traditionally a subfield of statistics) is machine learning (ML) which is traditionally a subfield of computer science. ML used to focus more on algorithms rather on probabilistic modeling but nowadays most machine learning methods are fully based on statistical multivariate approaches, so the two fields are converging.

Learning multivariate models allows us to learn dependencies and interactions among the components of the random variables which in turns allows to draw conclusion about the world.

Two main tasks:

- unsupervised learning (finding structure, clustering)
- supervised learning (training from labeled data, followed by prediction)

Challenges:

- complexity of model needs to be appropriate for problem and available data,
- high dimensions make estimation and inference difficult
- computational issues.

1.2 Basics

1.2.1 Univariate vs. multivariate random variables

Univariate random variable (dimension $d = 1$):

$$x \sim F$$

where x is a **scalar** and F is the distribution. $E(x) = \mu$ denotes the mean and $\text{Var}(x) = \sigma^2$ the variance of x .

Multivariate random **vector** of dimension d :

$$\mathbf{x} = (x_1, x_2, \dots, x_d)^T \sim F$$

\mathbf{x} is **vector** valued random variable.

The vector \mathbf{x} is column vector (=matrix of size $d \times 1$). Its components x_1, x_2, \dots, x_d are univariate random variables. The dimension d is also often denoted by p or q .

1.2.2 Mean of a random vector

The mean / expectation of a random vector with dimensions d is also a vector with dimensions d :

$$E(\mathbf{x}) = \boldsymbol{\mu} = \begin{pmatrix} E(x_1) \\ E(x_2) \\ \vdots \\ E(x_d) \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{pmatrix}$$

1.2.3 Variance of a random vector

Recall the definition of variance for a univariate random variable:

$$\text{Var}(x) = E((x - E(x))^2) = E((x - \mu)^2) = E((x - \mu)(x - \mu)) = E(x^2) - \mu^2$$

Definition of **variance of a random vector**:

$$\text{Var}(\mathbf{x}) = E \left(\underbrace{(\mathbf{x} - \boldsymbol{\mu})}_{d \times 1} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^T}_{1 \times d} \right) = \underbrace{\boldsymbol{\Sigma}}_{d \times d} = E(\mathbf{x}\mathbf{x}^T) - \boldsymbol{\mu}\boldsymbol{\mu}^T$$

The variance of a random vector is, therefore, **not** a vector but a **matrix**!

$$\boldsymbol{\Sigma} = (\sigma_{ij}) = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_{dd} \end{pmatrix}$$

This matrix is called the **Covariance Matrix**, with off-diagonal elements $\sigma_{ij} = \text{Cov}(x_i, x_j)$ and the diagonal $\sigma_{ii} = \text{Var}(X_i) = \sigma_i^2$.

1.2.4 Properties of the covariance matrix

1. Σ is real valued: $\sigma_{ij} \in \mathbb{R}$
2. Σ is symmetric: $\sigma_{ij} = \sigma_{ji}$
3. The diagonal of Σ contains $\sigma_{ii} = \text{Var}(x_i) = \sigma_i^2$, i.e. the variances of the components of x .
4. Off-diagonal elements $\sigma_{ij} = \text{Cov}(x_i, x_j)$ represent linear dependencies among the x_i . \implies linear regression, correlation

How many entries does the Σ matrix have?

$$\Sigma = (\sigma_{ij}) = \underbrace{\begin{pmatrix} \sigma_{11} & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_{dd} \end{pmatrix}}_{d \times d}$$

Number of entries: $\frac{d(d+1)}{2}$, which grows with square of dimension d , i.e. is of order $O(d^2)$.

d	# entries
1	1
10	55
100	5050
1000	500500
10000	50005000

For large dimension d the covariance matrix has many components!

\rightarrow computationally expensive (both for storage and in handling) \rightarrow very challenging to estimate in high dimensions d .

Note: matrix inversion requires $O(d^3)$ operations! So, computing Σ^{-1} is difficult for large d !

1.2.5 Eigenvalue decomposition of Σ

Theorem from matrix theory / linear algebra: A symmetric matrix with real entries has real eigenvalues.

Thus, the eigenvalues of Σ must be real. However, for covariance matrices this can be clarified further:

Assume that $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ is the eigenvalue decomposition of the covariance matrix Σ with \mathbf{U} an orthogonal matrix containing the eigenvectors (eigensystem) and $\mathbf{\Lambda}$ is the diagonal matrix containing eigenvalues:

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_d \end{pmatrix}$$

Assume that x is a random vector with covariance $\text{Var}(x) = \Sigma$. Then $\text{Var}(U^T x) = U^T \Sigma U = \Lambda$. Since the variance is always positive or zero the eigenvalues λ_i of the covariance matrix Σ cannot be negative. Hence, Σ is **positive semidefinite**.

In fact, **unless there is collinearity** (i.e. a variable is a linear function the other variables) all eigenvalues will be positive and Σ is **positive definite**.

1.2.6 Quantities related to the covariance matrix

1.2.6.1 Correlation matrix P

The correlation matrix P (= upper case greek “rho”) is the standardised covariance matrix

$$\rho_{ij} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}\sigma_{jj}}} = \text{Cor}(x_i, x_j)$$

$$\rho_{ii} = 1 = \text{Cor}(x_i, x_i)$$

$$P = (\rho_{ij}) = \begin{pmatrix} 1 & \dots & \rho_{1d} \\ \vdots & \ddots & \vdots \\ \rho_{d1} & \dots & 1 \end{pmatrix}$$

where P (“capital rho”) is a symmetric matrix ($\rho_{ij} = \rho_{ji}$).

Note the **variance-correlation decomposition**

$$\Sigma = V^{\frac{1}{2}} P V^{\frac{1}{2}}$$

where V is a diagonal matrix containing the variances:

$$V = \begin{pmatrix} \sigma_{11} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{dd} \end{pmatrix}$$

$$P = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}}$$

This is the definition of correlation written in matrix notation.

1.2.6.2 Precision matrix or concentration matrix

$$\Omega = (\omega_{ij}) = \Sigma^{-1}$$

Ω (“Omega”) is the inverse of the covariance matrix.

The inverse of the covariance matrix can be obtained via the spectral decomposition, followed by inverting the eigenvalues λ_i :

$$\Sigma^{-1} = \mathbf{U}\Lambda^{-1}\mathbf{U}^T = \mathbf{U} \begin{pmatrix} \lambda_1^{-1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_d^{-1} \end{pmatrix} \mathbf{U}^T$$

Note that **all eigenvalues λ_i need to be positive so that Σ can be inverted.** (i.e., Σ needs to be positive definite).

If any $\lambda_i = 0$ then Σ is singular and not invertible.

Importance of Σ^{-1} : - Natural parameter in exponential family. - Many expressions in multivariate statistics contain Σ^{-1} and not Σ - Σ^{-1} has close connection with graphical models (e.g. conditional independence graph, partial correlations, see later chapter)

1.2.6.3 Partial correlation matrix

This is a standardised version of the precision matrix, see later chapter on graphical models.

1.2.6.4 Total variation and generalised variance

To summarise the covariance matrix Σ in a single scalar value there are two commonly used measures:

- **total variation:** $\text{Tr}(\Sigma) = \sum_{i=1}^d \lambda_i$
- **generalised variance:** $\det(\Sigma) = \prod_{i=1}^d \lambda_i$

If all eigenvalues are positive then $\log \det(\Sigma) = \sum_{i=1}^d \log \lambda_i = \text{Tr}(\log \Sigma)$.

$\log \Sigma$ is the matrix logarithm of Σ and is given by $\log \Sigma = \mathbf{U} \begin{pmatrix} \log \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \log \lambda_d \end{pmatrix} \mathbf{U}^T$

1.3 Multivariate normal distribution

The multivariate normal model is a generalisation of the univariate normal distribution from dimension 1 to dimension d .

1.3.1 Univariate normal distribution:

Dimension $d = 1$

$$x \sim N(\mu, \sigma^2)$$

$$\mathbb{E}(x) = \mu, \text{Var}(x) = \sigma^2$$

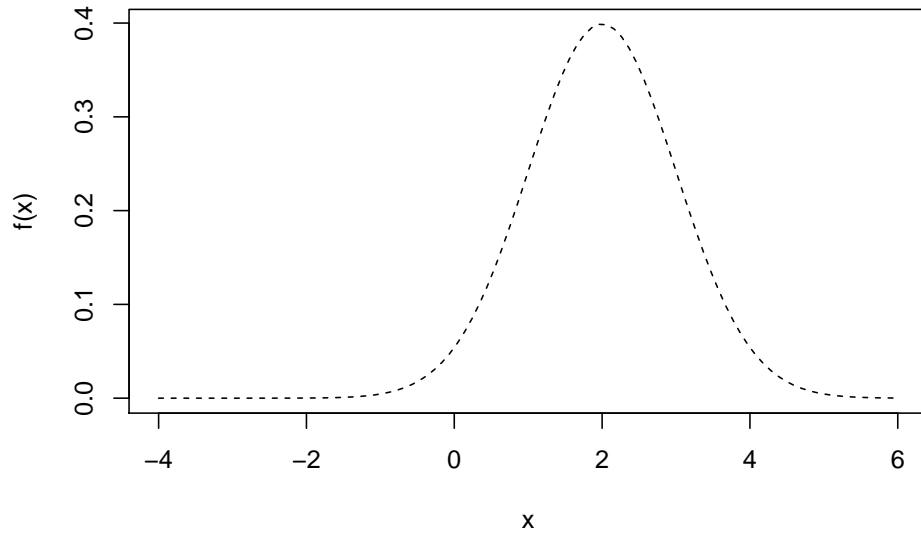
Density:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Plot of univariate normal density :

- Unimodal with peak at μ , the width is determined by σ (in this plot: $\mu = 2, \sigma = 1$)

Density Normal Distribution



Special case: **standard normal** with $\mu = 0$ and $\sigma^2 = 1$:

$$f(x|\mu = 0, \sigma^2 = 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

Maximum entropy characterisation: the normal distribution is the unique distribution that has the highest (differential) entropy over all continuous distributions with support from minus infinity to plus infinity with a given mean and variance.

This is in fact one of the reasons why the normal distribution is so important (and useful) – if we only know that a random variable has a mean and variance, and not much else, then using the normal distribution will be a reasonable and well justified working assumption!

1.3.2 Multivariate normal model

Dimension d

$$x \sim N_d(\mu, \Sigma)$$

$$x \sim \text{MVN}(\mu, \Sigma)$$

$$E(x) = \mu, \text{Var}(x) = \Sigma$$

Density:

$$f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp \left(-\frac{1}{2} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^T}_{1 \times d} \underbrace{\boldsymbol{\Sigma}^{-1}}_{d \times d} \underbrace{(\mathbf{x} - \boldsymbol{\mu})}_{d \times 1} \right)$$

$1 \times 1 = \text{scalar!}$

- note that density contains precision matrix $\boldsymbol{\Sigma}^{-1}$
- inverting $\boldsymbol{\Sigma}$ implies inverting the eigenvalues λ_i of $\boldsymbol{\Sigma}$ (thus we need $\lambda_i > 0$)
- density also contains $\det(\boldsymbol{\Sigma}) = \prod_{i=1}^d \lambda_i \equiv$ product of eigenvalues of $\boldsymbol{\Sigma}$

Special case: **standard multivariate normal** with

$$\boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma} = \mathbf{I} = \begin{pmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{pmatrix}$$

$$f(\mathbf{x}|\boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma} = \mathbf{I}) = (2\pi)^{-d/2} \exp \left(-\frac{1}{2} \mathbf{x}^T \mathbf{x} \right) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{x_i^2}{2} \right)$$

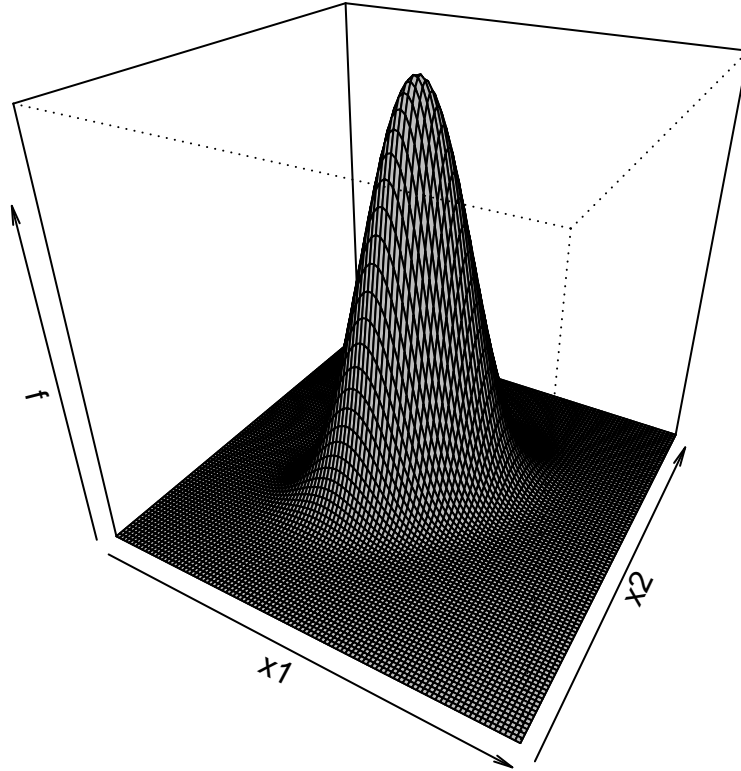
which is equivalent to the product of d univariate standard normals!

Misc:

- for $d = 1$, multivariate normal reduces to normal.
- for $\boldsymbol{\Sigma}$ diagonal (i.e. $\mathbf{P} = \mathbf{I}$, no correlation), MVN is the product of univariate normals (see Worksheet~2).

Plot of MVN density:

Density of bivariate normal (d=2)



- Location: μ
- Shape: Σ
- Unimodal: **one** peak
- Support from $-\infty$ to $+\infty$ in each dimension

1.3.3 Shape of the contour lines of the multivariate normal density

Now show that the contour lines of the multivariate normal density always take on the form of an ellipse, and that the radii of the ellipse is determined by the eigenvalues of Σ .

We start by observing that a circle with radius r around the origin can be described as the set of points (x_1, x_2) satisfying $x_1^2 + x_2^2 = r^2$, or equivalently, $\frac{x_1^2}{r^2} + \frac{x_2^2}{r^2} = 1$. This is generalised to the shape of an ellipse by allowing (in two dimensions) for two radii r_1 and r_2 with $\frac{x_1^2}{r_1^2} + \frac{x_2^2}{r_2^2} = 1$, or in vector notation

$\mathbf{x}^T \text{Diag}(r_1^2, r_2^2)^{-1} \mathbf{x} = 1$. In d dimensions and allowing for rotation of the axes and a shift of the origin from 0 to $\boldsymbol{\mu}$ the condition for an ellipse is

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q} \text{Diag}(r_1^2, \dots, r_d^2)^{-1} \mathbf{Q}^T (\mathbf{x} - \boldsymbol{\mu}) = 1$$

where \mathbf{Q} is an orthogonal rotation matrix.

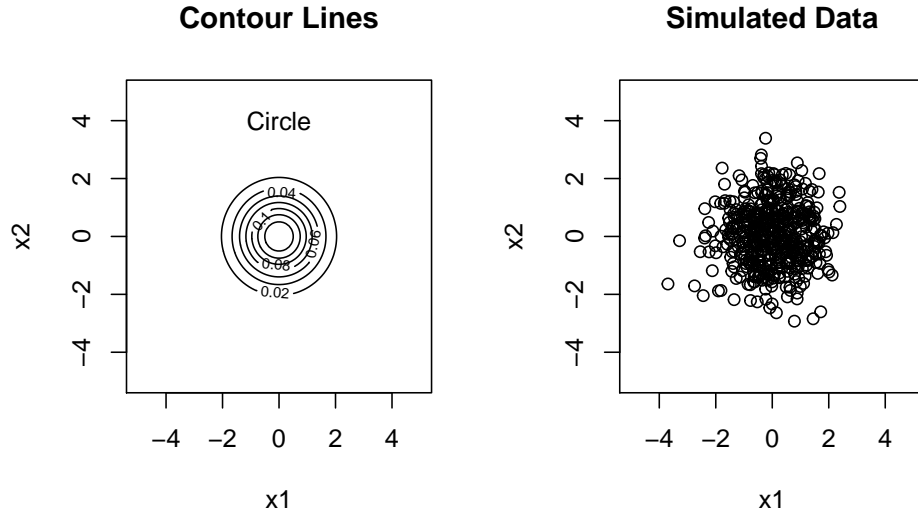
A contour line of a probability density function is a set of connected points where the density assumes the same constant value. In the case of the multivariate normal distribution keeping the density at some fixed value implies that $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = c$ where c is a constant. Using the eigenvalue decomposition of $\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$ we can rewrite this condition as

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu}) = c.$$

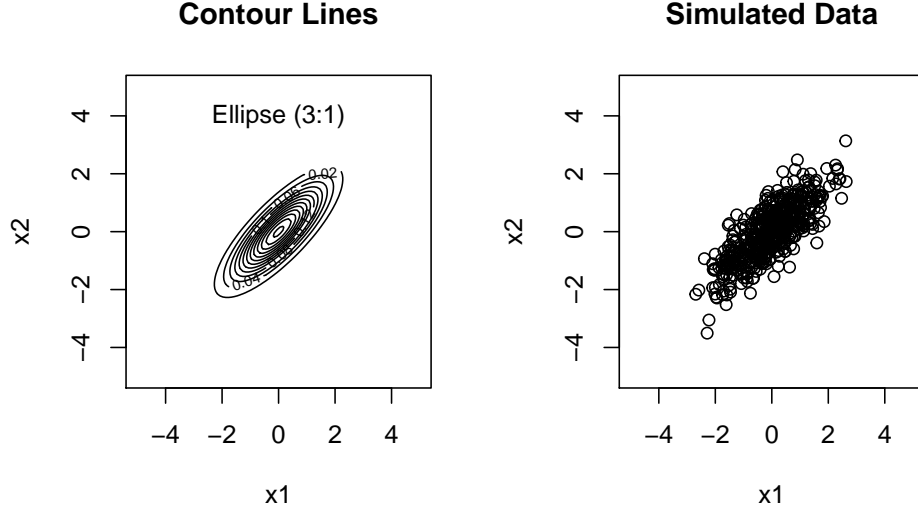
This implies that contour lines of the multivariate normal density are indeed ellipses and that the squared radii are proportional to the eigenvalues of $\boldsymbol{\Sigma}$. Equivalently, the positive square roots of the eigenvalues are proportional to the radii of the ellipse. Hence, for singular covariance matrix with one or more $\lambda_i = 0$ the corresponding radii are zero.

Two examples:

Case 1: No Correlation / Diagonal / Isotropic / Spherical Covariance $\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ with $\sqrt{\lambda_1/\lambda_2} = 1$:



Case 2: with correlation $\Sigma = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$ with $\sqrt{\lambda_1/\lambda_2} = 3$:



Small λ_i indicates collinearity!

1.4 Estimation in large sample and small sample settings

In practical application of multivariate normal model we need to learn its parameters from data. We first consider the case when there are many measurements available, and then second the case when the number of data points is small compared to the number of parameters.

1.4.1 Data matrix

Observations from a multivariate normal are vectors:

$$x_1, x_2, \dots, x_n \stackrel{\text{iid}}{\sim} N_d(\mu, \Sigma)$$

Data matrix (statistics convention):

Each *line* of the matrix is a transposed vector x_k^T .

Thus:

$$X = (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & & & \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{pmatrix}$$

with

$$x_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1d} \end{pmatrix}, x_2 = \begin{pmatrix} x_{21} \\ \vdots \\ x_{2d} \end{pmatrix}, \dots, x_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{nd} \end{pmatrix}$$

1.4. ESTIMATION IN LARGE SAMPLE AND SMALL SAMPLE SETTINGS¹⁷

Thus, in statistics the first index runs over $(1, \dots, n)$ and denotes the samples while the second index runs over $(1, \dots, d)$ and refers to the variables.

The convention on data matrices that variables are in columns while samples are in rows. is *not* universal! In fact it's the **other way around in machine learning** (where samples are stored in columns and variables in rows). However, some machine learning books also follow the statistics convention.

1.4.2 Strategies for large sample estimation

1.4.2.1 Empirical estimators (outline)

For large n :

$$\underbrace{F}_{\text{true}} \approx \underbrace{\hat{F}}_{\text{empirical}}$$

Replacing F by \hat{F} leads to *empirical estimators*.

For example, the expectation can be approximated/estimated as follows:

$$E_F(\mathbf{x}) \approx E_{\hat{F}}(\mathbf{x}) = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$E_F(g(\mathbf{x})) \approx E_{\hat{F}}(g(\mathbf{x})) = \frac{1}{n} \sum_{k=1}^n g(\mathbf{x}_k)$$

Recipe: to obtain an empirical plug-in estimator simply replace the expectation by the sample average in the population expression of the quantity of interest.

What does this work: the empirical distribution \hat{F} is actually the nonparametric maximum likelihood estimate of F (see below for likelihood estimation).

Note: the approximation of F by \hat{F} also the basis other approaches such as Efron's bootstrap method.

1.4.2.2 Maximum likelihood estimation (outline)

R.A. Fisher (1922): model-based estimators using the density or probability mass function

log-likelihood function:

$$\log L(\theta) = \sum_{k=1}^n \underbrace{f}_{\text{density}} \left(\underbrace{x_i}_{\text{data}} \mid \underbrace{\theta}_{\text{parameters}} \right)$$

= conditional probability of the observed data given the model parameters

Maximum likelihood estimate:

$$\hat{\theta}^{\text{ML}} = \arg \max_{\theta} \log L(\theta)$$

ML finds the parameters that make the observed data most likely (it does *not* find the most probable model!)

The great appeal of **MLEs** is that they **are optimal for large n**, i.e. they use all the information available in the data optimally to estimate parameters, and **for large sample size no estimator can be constructed that outperforms the MLE!**

A further advantage of the method of maximum likelihood is that it does not only provide a point estimate but also the asymptotic error (via the Fisher information which is related to the curvature of the log-likelihood function).

1.4.3 Large sample estimates of mean μ and covariance Σ

1.4.3.1 Empirical estimates:

These can be written in three different ways:

Vector notation

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\underbrace{\hat{\Sigma}}_{d \times d} = \frac{1}{n} \sum_{k=1}^n \underbrace{(x_k - \hat{\mu})}_{d \times 1} \underbrace{(x_k - \hat{\mu})^T}_{1 \times d}$$

Component notation

$$\hat{\mu}_i = \frac{1}{n} \sum_{k=1}^n x_{ki}$$

$$\hat{\sigma}_{ij} = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \hat{\mu}_i) (x_{kj} - \hat{\mu}_j)$$

$$\hat{\mu} = \begin{pmatrix} \hat{\mu}_1 \\ \vdots \\ \hat{\mu}_d \end{pmatrix}, \hat{\Sigma} = (\hat{\sigma}_{ij})$$

Variance estimate:

$$\hat{\sigma}_{ii} = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \hat{\mu}_i)^2$$

Note the factor $\frac{1}{n}$ (not $\frac{1}{n-1}$)

Data matrix notation

The empirical mean and covariance can also be written in terms of the data matrix X . See Worksheet~2 for details.

1.4.3.2 Maximum likelihood estimates

It turns out that the empirical estimates are identical to the MLE assuming multivariate normal model:

\implies empirical estimates $\hat{\mu}$ and $\hat{\Sigma}$

$$\hat{\mu}^{\text{ML}} = \hat{\mu}^{\text{emp}}$$

$$\hat{\Sigma}^{\text{ML}} = \hat{\Sigma}^{\text{emp}}$$

For a direct derivation of the multivariate normal MLEs by optimising the multivariate normal log-likelihood function see the Worksheet~2 (easy for the mean, more difficult for the covariance matrix).

Note the factor $\frac{1}{n}$ in the MLE of the covariance matrix.

1.4.3.3 Distribution of the empirical / maximum likelihood estimates

With $x_1, \dots, x_n \sim N_d(\mu, \Sigma)$ one can find the exact distributions of the estimators.

1. Distribution of the estimate of the mean:

$$\hat{\mu} \sim N_d\left(\mu, \frac{\Sigma}{n}\right)$$

Since $E(\hat{\mu}) = \mu \implies \hat{\mu}$ is unbiased

2. Distribution of the covariance estimate:

$$n\hat{\Sigma} \sim \text{Wishart}(\Sigma, n-1)$$

Since $E(n\hat{\Sigma}) = (n-1)\Sigma \implies \hat{\Sigma}$ is biased!

Easy to make unbiased: $\frac{n}{n-1}\hat{\Sigma} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^T$ is unbiased

But unbiasedness is **not** a very relevant criteria in multivariate statistics!

1.4.4 Problems with maximum likelihood in small sample settings and high dimensions

Modern data is high dimensional!

Data sets with $n < d$, i.e. high dimension d and small sample size n are now common in many fields, e.g., medicine, biology but also finance and business analytics.

$$n = 100 \text{ (e.g. patients/samples)}$$

$$d = 20000 \text{ (e.g., genes/SNPs/proteins/variables)}$$

Reasons:

- the number of measured variables is increasing quickly with technological advances (e.g. genomics)

- but the number of samples cannot be similarly increased (for cost and ethical reasons)

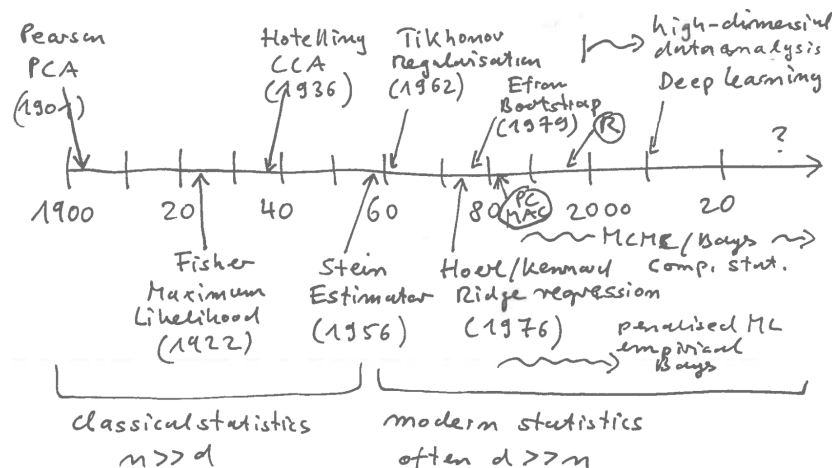
General problems of MLEs:

1. ML estimators are optimal only if **sample size is large** compared to the number of parameters. However, this optimality is not any more valid if sample size is moderate or smaller than the number of parameters.
2. If there is not enough data the **ML estimate overfits**. This means ML fits the current data perfectly but the resulting model does not generalise well (i.e. model will perform poorly in prediction)
3. If there is a choice between different models with different complexity **ML will always select the model with the largest number of parameters**.

-> for high-dimensional data with small sample size maximum likelihood estimation does not work!!!

History of Statistics:

Much of modern statistics (from 1960 onwards) is devoted to the development of inference and estimation techniques that work with complex, high-dimensional data.



- Maximum likelihood is a method from classical statistics (time up to about 1960).
- From 1960 modern (computational) statistics emerges, starting with “**Stein Paradox**” (1956): Charles Stein showed that in a **multivariate setting** ML estimators are **dominated by** (= are always worse than) shrinkage estimators!
- For example, there is a shrinkage estimator for the mean that is better (in terms of MSE) than the average (which is the MLE)!

Modern statistics has developed many different (but related) methods for use in high-dimensional, small sample settings:

- regularised estimators
- shrinkage estimators
- penalised maximum likelihood estimators

1.4. ESTIMATION IN LARGE SAMPLE AND SMALL SAMPLE SETTINGS 21

- Bayesian estimators
- Empirical Bayes estimators
- KL / entropy-based estimators

Most of this is out of scope for our class, but will be covered in advanced statistical courses.

Next, we describe a **simple regularised estimator for the estimation of the covariance** that we will use later (i.e. in classification).

1.4.5 Estimation of covariance matrix in small sample settings

Problems with ML estimate of Σ

1. Σ has $O(d^2)$ number of parameters! $\implies \hat{\Sigma}^{\text{MLE}}$ requires *a lot* of data!
 $n \gg d$ or d^2
2. if $n < d$ then $\hat{\Sigma}$ is positive **semi**-definite (even if Σ is p.d.f.!)
 $\implies \hat{\Sigma}$ will have **vanishing eigenvalues** (some $\lambda_i = 0$) and thus **cannot be inverted** and is singular!

Simple regularised estimate of Σ

Regularised estimator S^* = convex combination of $S = \hat{\Sigma}^{\text{MLE}}$ and I_d (identity matrix) to get

Regularisation:

$$\underbrace{S^*}_{\text{regularised estimate}} = \underbrace{\lambda}_{\text{shrinkage intensity}} \underbrace{I_d}_{\text{target}} + (1 - \lambda) \underbrace{S}_{\text{ML estimate}}$$

Next, choose $\lambda \in [0, 1]$ such that S^* is better (in terms of MSE) than both S and I_d .

Bias-variance trade-off

MSE is the Mean Squared Error, composed of squared bias and variance.

$$\text{MSE}(\theta) = E((\hat{\theta} - \theta)^2) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$$

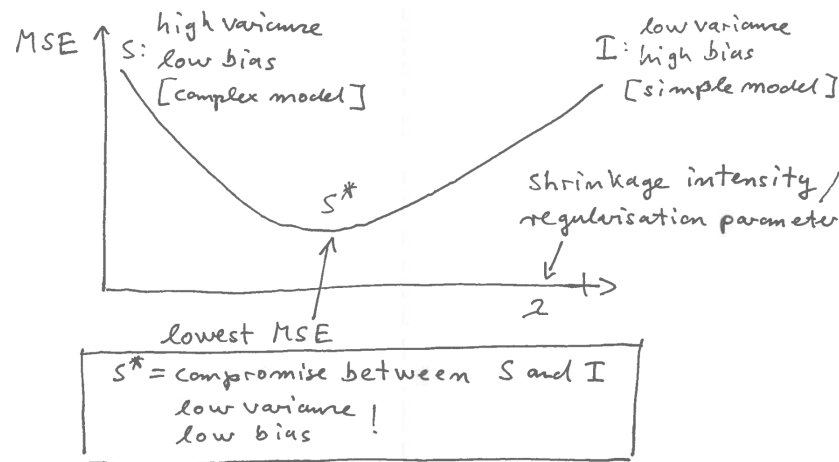
with $\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$

S : ML estimate, many parameters, low bias, high variance

I_d : “target”, no parameters, high bias, low variance

\implies **reduce high variance of S by introducing a bit of bias through I_d !**

\implies overall, MSE is decreased



How to find optimal shrinkage / regularisation parameter λ ? Minimise MSE!

Challenge: since we don't know the true Σ we cannot actually compute the MSE directly but have to estimate it! How is this done in practise?

- by cross-validation (=resampling procedure)
- by using some analytic approximation (e.g. Stein's formula)

Why does regularisation of $\hat{\Sigma}$ work?

1. S^* is **positive definite**:

Matrix Theory:

$$\underbrace{M_1}_{\text{positive definite, } \lambda I} + \underbrace{M_2}_{\text{positive semi-definite, } (1-\lambda)S} = \underbrace{M_3}_{\text{positive definite, } S^*}$$

$\Rightarrow S^*$ can be inverted even if $n < d$

2. It's **Bayesian** in disguise!

$$\underbrace{S^*}_{\text{posterior mean}} = \underbrace{\lambda I_d}_{\text{prior information}} + (1-\lambda) \underbrace{S}_{\text{data summarised by maximum likelihood}}$$

- Prior information helps to infer Σ even in small samples
- Since λ is chosen from data, it is actually an empirical Bayes.
- also called shrinkage estimator since the off-diagonal entries are shrunk towards zero
- this type of linear shrinkage/regularisation is natural for models in the exponential family (Diaconis and Ylvisaker, 1979)

In Worksheet~2 the empirical estimator of covariance is compared with the covariance estimator implemented in the R package "corpcor". This uses a regularisation similar as above (but for the correlation rather than covariance matrix) and it employs an analytic data-adaptive estimate of the shrinkage intensity λ . This estimator is a variant of an empirical Bayes / James-Stein estimator (see the year 2 Statistical Methods 20802 module).

Summary

- In multivariate statistics, it is useful (and often necessary) to utilise prior information!
- Regularisation introduces bias and reduces variance, minimising overall MSE
- Unbiased estimation (a highly valued property in classical statistics!) is not a good idea in multivariate settings and often leads to poor estimators.

1.5 Discrete multivariate distributions

Most univariate distributions have multivariate counterparts. Here are some of the most important ones. First we discuss discrete distributions, then later further continuous distributions

1.5.1 Categorical distribution

1.5.1.1 Univariate case

Bernoulli distribution (=coin tossing model)

$$\begin{aligned}
 x &\sim \text{Ber}(\pi) \\
 x &\in \{0, 1\} \\
 \pi &\in [0, 1] \\
 \Pr(x = 1) &= \pi \\
 \Pr(x = 0) &= 1 - \pi \\
 E(x) &= \pi \\
 \text{Var}(x) &= \pi(1 - \pi)
 \end{aligned}$$

1.5.1.2 Multivariate case

$$\begin{aligned}
 \mathbf{x} &\sim \text{Categ}(\boldsymbol{\pi}) \\
 \mathbf{x} &= (x_1, \dots, x_d)^T; x_i \in \{0, 1\}; \sum_{i=1}^d x_i = 1 \\
 \boldsymbol{\pi} &= (\pi_1, \dots, \pi_d)^T; \sum_{i=1}^d \pi_i = 1 \\
 \Pr(x_i = 1) &= \pi_i \\
 \Pr(x_i = 0) &= 1 - \pi_i \\
 E(\mathbf{x}) &= \boldsymbol{\pi} \\
 \text{Var}(x_i) &= \pi_i(1 - \pi_i) \\
 \text{Cov}(x_i, x_j) &= -\pi_i\pi_j
 \end{aligned}$$

1.5.2 Multinomial distribution

1.5.2.1 Univariate case

Binomial Distribution

Repeat Bernoulli experiment r times

$$x \sim \text{Binom}(\pi, r)$$

$$x \in \{0, \dots, r\}$$

$$E(x) = r \pi$$

$$\text{Var}(x) = r \pi(1 - \pi)$$

Standardised to unit interval:

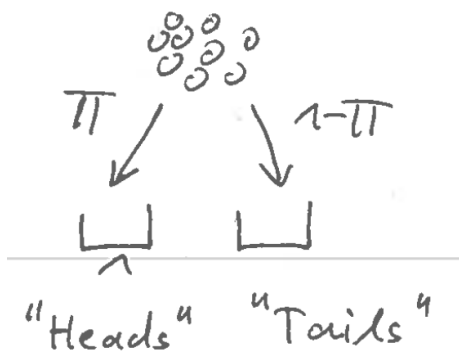
$$\frac{x}{r} \in \left\{0, \frac{1}{r}, \dots, 1\right\}$$

$$E\left(\frac{x}{r}\right) = \pi$$

$$\text{Var}\left(\frac{x}{r}\right) = \frac{\pi(1 - \pi)}{r}$$

Urn model:

distribute r balls into two bins



1.5.2.2 Multivariate case

Multinomial distribution

Draw r times from categorical distribution

$$x \sim \text{Multinom}(\pi, r)$$

$$x_i \in \{0, 1, \dots, r\}; \sum_{i=1}^d x_i = r$$

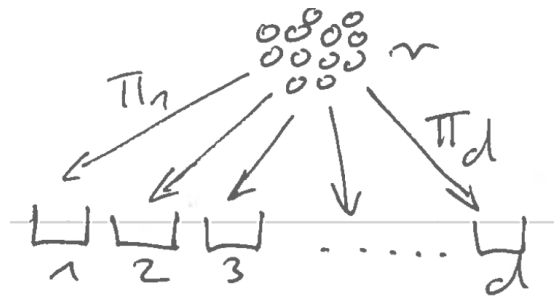
$$\begin{aligned} E(\mathbf{x}) &= r \boldsymbol{\pi} \\ \text{Var}(x_i) &= r \pi_i (1 - \pi_i) \\ \text{Cov}(x_i, x_j) &= -r \pi_i \pi_j \end{aligned}$$

Standardised to unit interval:

$$\begin{aligned} \frac{x_i}{r} &\in \left\{ 0, \frac{1}{r}, \frac{2}{r}, \dots, 1 \right\} \\ E\left(\frac{\mathbf{x}}{r}\right) &= \boldsymbol{\pi} \\ \text{Var}\left(\frac{x_i}{r}\right) &= \frac{\pi_i(1 - \pi_i)}{r} \\ \text{Cov}\left(\frac{x_i}{r}, \frac{x_j}{r}\right) &= -\frac{\pi_i \pi_j}{r} \end{aligned}$$

Urn model:

distribute r balls into d bins



1.6 Continuous multivariate distributions

1.6.1 Dirichlet distribution

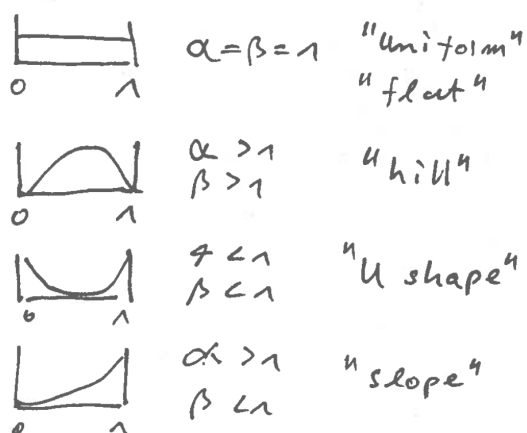
1.6.1.1 Univariate case

Beta distribution

$$\begin{aligned} x &\sim \text{Beta}(\alpha, \beta) \\ x &\in [0, 1] \\ \alpha &> 0; \beta > 0 \\ m &= \alpha + \beta \\ \mu &= \frac{\alpha}{m} \in [0, 1] \\ E(x) &= \mu \\ \text{Var}(x) &= \frac{\mu(1 - \mu)}{m + 1} \end{aligned}$$

compare with unit standardised binomial!

Different shapes



Useful as distribution for a proportion π

Bayesian Model:

Beta prior: $\pi \sim \text{Beta}(\alpha, \beta)$

Binomial likelihood: $x|\pi \sim \text{Binom}$

1.6.1.2 Multivariate case

Dirichlet distribution

$$\mathbf{x} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$x_i \in [0, 1]; \sum_{i=1}^d x_i = 1$$

$$\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^T > 0$$

$$m = \sum_{i=1}^d \alpha_i$$

$$\mu_i = \frac{\alpha_i}{m} \in [0, 1]$$

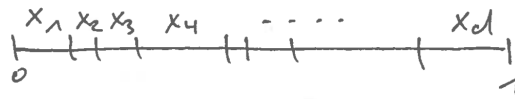
$$\mathbb{E}(x_i) = \mu_i$$

$$\text{Var}(x_i) = \frac{\mu_i(1 - \mu_i)}{m + 1}$$

$$\text{Cov}(x_i, x_j) = -\frac{\mu_i \mu_j}{m + 1}$$

compare with unit standardised multinomial!

Stick breaking" model



Useful as distribution for a proportion π

Bayesian Model:

Dirichlet prior: $\pi \sim \text{Dirichlet}(\alpha)$

Multinomial likelihood: $x|\pi \sim \text{Multinom}$

1.6.2 Wishart distribution

This is a distribution for the sum of squared normally distributed random variables.

1.6.2.1 Univariate case

Scaled χ^2 distribution

$$z_1, z_2, \dots, z_m \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$x = \sum_{i=1}^m z_i^2$$

$$x \sim \sigma^2 \chi_m^2 = W_1(\sigma^2, m)$$

$$E(x) = m \sigma^2$$

$$\text{Var}(x) = m 2 \sigma^4$$

Useful as the distribution of sample variance:

$$y_1, \dots, y_n \sim N(\mu, \sigma^2)$$

Known mean μ

$$\frac{1}{k} \sum_{i=1}^n (y_i - \mu)^2 \sim W_1\left(\frac{\sigma^2}{k}, n\right)$$

Unknown mean μ

$$\frac{1}{k} \sum_{i=1}^n (y_i - \bar{y})^2 \sim W_1\left(\frac{\sigma^2}{k}, n-1\right)$$

1.6.2.2 Multivariate case

Wishart distribution

$$z_1, z_2, \dots, z_m \stackrel{\text{iid}}{\sim} N_d(0, \Sigma)$$

$$\underbrace{\mathbf{X}}_{d \times d} = \sum_{i=1}^m \underbrace{z_i z_i^T}_{d \times d}$$

Note that \mathbf{X} is a *matrix*!

$$\mathbf{X} \sim W_d(\Sigma, m)$$

$$E(\mathbf{X}) = m\Sigma$$

$$\text{Var}(x_{ij}) = m \left(\sigma_{ij}^2 + \sigma_{ii}\sigma_{jj} \right)$$

Useful as distribution of sample covariance:

$$\mathbf{y}_1, \dots, \mathbf{y}_n \sim N_d(\boldsymbol{\mu}, \Sigma)$$

$$\frac{1}{k} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T \sim W_d(\Sigma/k, n)$$

$$\frac{1}{k} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T \sim W_d(\Sigma/k, n-1)$$

1.6.2.3 Relationship to Gamma distribution

The scaled χ^2 distribution (=one-dimensional Wishart distribution) with parameters σ^2 and m is in fact a reparameterised **Gamma distribution** with shape parameter α and scale parameter β :

$$\text{Gamma} \left(\underbrace{\frac{m}{2}}_{\text{shape}}, \underbrace{2\sigma^2}_{\text{scale}} \right) = \sigma^2 \chi_m^2 = W_1(\sigma^2, m)$$

or, equivalently ($m = 2\alpha, \sigma^2 = \beta/2$)

$$\text{Gamma} \left(\underbrace{\alpha}_{\text{shape}}, \underbrace{\beta}_{\text{scale}} \right) = \frac{\beta}{2} \chi_{2\alpha}^2 = W_1\left(\frac{\beta}{2}, 2\alpha\right)$$

The mean of the Gamma distribution is $E(x) = \alpha\beta = \mu$ and the variance is $\text{Var}(x) = \alpha\beta^2 = \mu^2/\alpha$.

The **exponential distribution** with rate parameter λ is a special case of the Gamma distribution with $\alpha = 1$:

$$\text{Exp}(\lambda) = \text{Gamma}(1, \frac{1}{\lambda}) = \frac{1}{2\lambda} \chi^2_2 = W_1(\frac{1}{2\lambda}, 2)$$

The corresponding mean is $1/\lambda = \mu$ and variance $1/\lambda^2 = \mu^2$.

1.6.3 Inverse Wishart distribution

1.6.3.1 Univariate case

Inverse χ^2 Distribution

$$x \sim W_1^{-1}(\psi, k+2) = \psi \text{ Inv-}\chi^2_{k+2}$$

$$E(x) = \frac{\psi}{k}$$

$$\text{Var}(x) = \frac{2\psi^2}{k^2(k-2)}$$

Relationship to scaled χ^2 :

$$\frac{1}{x} \sim W_1(\psi^{-1}, k+2) = \psi^{-1} \chi^2_{k+2}$$

1.6.3.2 Multivariate case

Inverse Wishart distribution

$$\underbrace{\mathbf{X}}_{d \times d} \sim W_d^{-1} \left(\underbrace{\mathbf{\Psi}}_{d \times d}, k+d+1 \right)$$

$$E(\mathbf{X}) = \mathbf{\Psi}/k$$

$$\text{Var}(x_{ij}) = \frac{2}{k^2(k-2)} \frac{(k+2)\psi_{ij} + k\psi_{ii}\psi_{jj}}{2k+2}$$

Relationship to Wishart:

$$\mathbf{X}^{-1} \sim W_d \left(\mathbf{\Psi}^{-1}, k+d+1 \right)$$

1.6.3.3 Relationship to inverse Gamma distribution

Another way to express the univariate inverse Wishart distribution is via the **inverse Gamma distribution**:

$$IG(\underbrace{1 + \frac{k}{2}}_{\text{shape } \alpha}, \underbrace{\frac{\psi}{2}}_{\text{scale } \beta}) = \psi \text{ Inv-}\chi^2_{k+2} = W_1^{-1}(\psi, k+2)$$

or equivalently ($k = 2(\alpha - 1)$ and $\psi = 2\beta$)

$$IG(\alpha, \beta) = 2\beta \text{ Inv-}\chi^2_{2\alpha} = W_1^{-1}(2\beta, 2\alpha)$$

The mean of the inverse Gamma distribution is $E(x) = \frac{\beta}{\alpha-1} = \mu$ the variance $\text{Var}(x) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)} = \frac{2\mu^2}{k-2}$.

The inverse of x is Gamma distributed:

$$\frac{1}{x} \sim \text{Gamma}\left(1 + \frac{k}{2}, 2\psi^{-1}\right) = \text{Gamma}(\alpha, \beta^{-1})$$

The inverse Wishart distribution is useful as conjugate distribution for Bayesian modelling of the variance, with k the sample size parameter and $\Psi = k\Sigma$ (or $\psi = k\sigma^2$).

1.6.4 Further distributions

https://en.wikipedia.org/wiki/List_of_probability_distributions

Wikipedia is a quite good source for information on distributions!

Chapter 2

Transformations and dimension reduction

Motivation: In the following we study transformations of random vectors and their distributions. These transformation are very important since they either transform simple distributions into more complex distributions or allow to simplify complex models. In machine learning invertible mappings of transformations for probability distributions are known as “normalising flows” (these play a key role e.g. in neural networks).

2.1 Linear Transformations

2.1.1 Location-scale transformation

Also known as affine transformation.

$$y = \underbrace{a}_{\text{location parameter}} + \underbrace{B}_{\text{scale parameter}} x$$

$y : m \times 1$ random vector

$a : m \times 1$ vector, location parameter

$B : m \times d$ matrix, scale parameter, $m \geq 1$

$x : d \times 1$ random vector

$$\begin{aligned} E(x) = \mu & \implies E(y) = a + B\mu \\ \text{Var}(x) = \Sigma & \implies \text{Var}(y) = B\Sigma B^T \end{aligned}$$

Special cases/examples:

1. Univariate case ($d = 1, m = 1$)

- $E(y) = a + b\mu$
 - $\text{Var}(y) = b^2\sigma^2$
2. Sum of two random univariate variables
 $y = x_1 + x_2$, i.e. $a = 0$ and $B = (1, 1)$
- $E(x_1 + x_2) = \mu_1 + \mu_2$
 - $\text{Var}(x_1 + x_2) = (1, 1) \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \sigma_1^2 + \sigma_2^2 + 2\sigma_{12} = \text{Var}(x_1) + \text{Var}(x_2) + 2\text{Cov}(x_1, x_2)$
3. $y_1 = a_1 + b_1x_1$ and $y_2 = a_2 + b_2x_2$, i.e. $a = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$ and $B = \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix}$
- $E(y) = \begin{pmatrix} a_1 + b_1\mu_1 \\ a_2 + b_2\mu_2 \end{pmatrix}$
 - $\text{Var}(y) = \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix} \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix} \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix} = \begin{pmatrix} b_1^2\sigma_1^2 & b_1b_2\sigma_{12} \\ b_1b_2\sigma_{21} & b_2^2\sigma_2^2 \end{pmatrix}$
i.e. $\text{Cov}(a_1 + b_1x_1, a_2 + b_2x_2) = b_1b_2\text{Cov}(x_1, x_2)$

2.1.2 Invertible location-scale transformation

If $m = d$ and $\det(B) \neq 0$ then we get an **invertible** transformation:

$$y = a + Bx$$

$$x = B^{-1}(y - a)$$

Transformation of density: $x \sim F_x$ with density $f_x(x)$

$\implies y \sim F_y$ with density

$$f_y(y) = |\det(B)|^{-1} f_x(B^{-1}(y - a))$$

Example 2.1. coloring transformation¹ $y = \mu + \Sigma^{1/2}x$

$$a = \mu$$

$$B = \Sigma^{1/2}$$

where $\Sigma^{1/2} = U\Lambda^{1/2}U^T$ is the principal matrix square root obtained by eigenvalue decomposition of $\Sigma = U\Lambda U^T$. Note that $\Sigma^{1/2}$ is unique, symmetric, and positive definite (not just positive semi-definite because this is an invertible transformation).

$$E(x) = \mathbf{0} \text{ and } \text{Var}(x) = I_d$$

$$\implies E(y) = \mu \text{ and } \text{Var}(y) = \Sigma$$

¹this is just one particular example of a coloring transformation. There are in fact infinitely many, see the discussion on whitening transformations.

Assume \mathbf{x} is multivariate standard normal $\mathbf{x} \sim N_d(\mathbf{0}, \mathbf{I}_d)$ with density

$$f_{\mathbf{x}}(\mathbf{x}) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2}\mathbf{x}^T \mathbf{x}\right)$$

Then the density after applying this coloring transform is

$$\begin{aligned} f_{\mathbf{y}}(\mathbf{y}) &= |\det(\mathbf{\Sigma}^{1/2})|^{-1} (2\pi)^{-d/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1/2} \mathbf{\Sigma}^{-1/2} (\mathbf{y} - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-d/2} \det(\mathbf{\Sigma})^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})\right) \end{aligned}$$

$\implies \mathbf{y}$ has multivariate normal density!!

Application: e.g. random number generation: draw from $N_d(\mathbf{0}, \mathbf{I}_d)$ (easy!) then convert to multivariate normal by coloring transformation $\mathbf{x} \rightarrow \mathbf{y}$

Example 2.2. Mahalanobis transform $\mathbf{y} = \mathbf{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu})$

This is the **inverse** of the above coloring transformation

$$\begin{aligned} \mathbf{a} &= -\mathbf{\Sigma}^{-1/2} \boldsymbol{\mu} \\ \mathbf{B} &= \mathbf{\Sigma}^{-1/2} \\ \mathbb{E}(\mathbf{x}) &= \boldsymbol{\mu} \text{ and } \text{Var}(\mathbf{x}) = \mathbf{\Sigma} \\ \implies \mathbb{E}(\mathbf{y}) &= \mathbf{0} \text{ and } \text{Var}(\mathbf{y}) = \mathbf{I}_d \end{aligned}$$

Mahalanobis transformation performs three functions:

1. Centering $(-\boldsymbol{\mu})$
2. Standardisation $\text{Var}(y_i) = 1$
3. Decorrelation $\text{Cor}(y_i, y_j) = 0$ for $i \neq j$

Univariate case ($d = 1$)

$$y = \frac{x - \mu}{\sigma}$$

As there is no correlation for $d = 1$, it's standardisation + centering.

Mahalanobis transformation appears implicitly in many places in multivariate statistics, e.g. in the multivariate normal density.

2.2 Nonlinear transformations

2.2.1 General transformation

$$\mathbf{y} = \mathbf{h}(\mathbf{x})$$

with \mathbf{h} an arbitrary vector-valued function

- linear case: $\mathbf{h}(\mathbf{x}) = \mathbf{a} + \mathbf{B}\mathbf{x}$

$$\begin{aligned} E(\mathbf{y}) &=? \\ \text{Var}(\mathbf{y}) &=? \end{aligned}$$

For a general transformation $h(\mathbf{x})$ the mean and variance of the transformed variable cannot be easily or analytically calculated.

However, we can find a **linear approximation** and then compute the mean and variance. This is called the “Delta Method”.

2.2.2 Linearisation of $h(\mathbf{x})$

Taylor series approximation (first order) of $h(\mathbf{x})$ around $\boldsymbol{\mu}$:

$$h(\mathbf{x}) \approx h(\boldsymbol{\mu}) + \underbrace{J_h(\boldsymbol{\mu})}_{\text{Jacobi matrix}} (\mathbf{x} - \boldsymbol{\mu})$$

∇ , the nabla operator, is the row vector $(\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_d})$, which when applied to univariate h gives the gradient:

$$\nabla h(\mathbf{x}) = \left(\frac{\partial h}{\partial x_1}, \dots, \frac{\partial h}{\partial x_d} \right)$$

The **Jacobi matrix** is the **generalisation of gradient** if h is vector-valued:

$$J_h(\mathbf{x}) = \begin{pmatrix} \nabla h_1(\mathbf{x}) \\ \nabla h_2(\mathbf{x}) \\ \vdots \\ \nabla h_m(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \frac{\partial h_1}{\partial x_1} & \dots & \frac{\partial h_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_m}{\partial x_1} & \dots & \frac{\partial h_m}{\partial x_d} \end{pmatrix}$$

2.2.3 Delta method

$$\mathbf{y} = h(\mathbf{x}) \approx \mathbf{a} + \mathbf{B}\mathbf{x}$$

with $\mathbf{a} = h(\boldsymbol{\mu}) - \boldsymbol{\mu} J_h(\boldsymbol{\mu})$ and $\mathbf{B} = J_h(\boldsymbol{\mu})$

\implies

$$E(\mathbf{y}) \approx h(\boldsymbol{\mu})$$

$$\text{Var}(\mathbf{y}) \approx J_h(\boldsymbol{\mu}) \boldsymbol{\Sigma} J_h(\boldsymbol{\mu})^T$$

Special case: univariate Delta method

$$E(y) = h(\mu)$$

$$\text{Var}(y) = \sigma^2 h'(\mu)^2$$

2.2.4 Transformation of densities under general invertible transformation

Assume $h(x) = y$ is invertible: $h^{-1}(y) = x(y)$

$x \sim F_x$ with probability density function $f_x(x)$

The density $f_y(y)$ of the transformed random vector y is then given by

$$f_y(y) = |\det(J_x(y))| f_x(x(y))$$

where $J_x(y)$ is the Jacobi matrix of the inverse transformation

Special cases:

- Univariate version: $f_y(y) = \left| \frac{dx}{dy} \right| f_x(x(y))$
- Linear transformation $h(x) = a + Bx$, with $x(y) = B^{-1}(y - a)$ and $J_x(y) = B^{-1}$:

$$f_y(y) = |\det(B)|^{-1} f_x(B^{-1}(y - a))$$

2.2.5 Normalising flows

In machine learning (sequences of) invertible nonlinear transformations are known as “normalising flows”. They are used both in a generative way (building complex models from simple models) and also in a simplification and dimension reduction context.

In this module we will focus mostly on linear transformations as these underpin much of classical multivariate statistics, but it is important to keep in mind for later study the importance of nonlinear transformations —see, e.g. the review paper by Kobyzev et al. “Normalizing Flows: Introduction and Ideas”, available from <https://arxiv.org/abs/1908.09257>.

2.3 Whitening transformations

2.3.1 Overview

The *Mahalanobis* transform (also known as ZCA transform in machine learning) is a specific example of a **whitening transformation**. These constitute an important and widely used class of invertible location-scale transformations.

Terminology: whitening refers to the fact that after the transformation the covariance matrix is spherical, isotropic, white (I_d)

Whitening is **useful in preprocessing**, to **turn multivariate problems into simple univariate models** and some **reduce the dimension in an optimal way**.

In so-called latent variable models whitening procedures link observed and latent variables (which usually are uncorrelated and standardised random variables):

Whitening $\begin{matrix} \downarrow \\ x \\ \uparrow \\ z \end{matrix}$ Observed external variable (can be measured), typically correlated
 Unobserved "latent" variable internal, typically not correlated

2.3.2 General whitening transformation

x random vector $\sim F_x$ (not necessarily from multivariate normal)

$$\underbrace{z}_{d \times 1 \text{ vector}} = \underbrace{W}_{d \times d \text{ whitening matrix}} \underbrace{x}_{d \times 1 \text{ vector}}$$

Objective: choose W so that $\text{Var}(z) = I_d$

Note: we do not care about $E(z)$ since we can always centre!

For Mahalanobis/ZCA whitening we already know that $W^{\text{ZCA}} = \Sigma^{-1/2}$.

In general, W needs to satisfy a constraint:

$$\begin{aligned} \text{Var}(z) &= I_d \\ \implies \text{Var}(Wx) &= W \Sigma W^T = I_d \\ \implies W \Sigma W^T W &= W \end{aligned}$$

$$\implies \text{constraint on whitening matrix: } W^T W = \Sigma^{-1}$$

Clearly, the ZCA whitening matrix satisfies this constraint: $(W^{\text{ZCA}})^T W^{\text{ZCA}} = \Sigma^{-1/2} \Sigma^{-1/2} = \Sigma^{-1}$

2.3.3 General solution of whitening constraint (covariance-based)

$$W = Q_1 \Sigma^{-1/2}$$

where Q_1 is a rotation (orthogonal) matrix (with properties $Q^T Q = Q Q^T = I$ and $Q^T = Q^{-1}$).

Easy to see that this W satisfies the whitening constraint:

$$W^T W = \Sigma^{-1/2} \underbrace{Q_1^T Q_1}_I \Sigma^{-1/2} = \Sigma^{-1}$$

\implies instead of choosing W , we choose the rotation matrix Q_1 !

- it is now clear that there are infinitely many whitening procedures, because there are infinitely many rotations! This also means we need to find ways to choose/select among whitening procedures.
- for the Mahalanobis/ZCA transformation $Q_1^{\text{ZCA}} = I$
- **whitening** can be interpreted as **Mahalanobis transform** followed by **rotation**

2.3.4 Another solution (correlation-based)

Instead of working with Σ , we can express W also in terms of correlation matrix $P = (\rho_{ij})$.

$$W = Q_2 P^{-1/2} V^{-1/2}$$

where $V^{1/2}$ is the diagonal matrix containing the variances.

It is easy to verify that this W also satisfies the whitening constraint:

$$\begin{aligned} W^T W &= V^{-1/2} P^{-1/2} \underbrace{Q_2^T Q_2}_I P^{-1/2} V^{-1/2} \\ &= V^{-1/2} P^{-1} V^{-1/2} = \Sigma^{-1} \end{aligned}$$

- for Mahalanobis/ZCA transformation $Q_2^{ZCA} = \Sigma^{-1/2} V^{1/2} P^{1/2}$
- **Another interpretation of whitening:** first **standardising** ($V^{-1/2}$), then **decorrelation** ($P^{-1/2}$), followed by **rotation** (Q_2)

Both forms to write W are equally valid (and interchangeable).

Note that for the same W

$$Q_1 \neq Q_2 \text{ Two different rotation matrices!}$$

and also

$$\underbrace{\Sigma^{-1/2}}_{\text{Symmetric}} \neq \underbrace{P^{-1/2} V^{-1/2}}_{\text{Not Symmetric}}$$

even though

$$\Sigma^{-1/2} \Sigma^{-1/2} = \Sigma^{-1} = V^{-1/2} P^{-1/2} P^{-1/2} V^{-1/2}$$

2.3.5 Objective criteria for choosing among W using Q_1 or Q_2

1. **Cross-covariance** $\Phi = \Sigma_{z,x}$ between z and x

$$\begin{aligned} \Phi = \text{Cov}(z, x) &= \text{Cov}(Wx, x) = W\Sigma \\ &= Q_1 \Sigma^{-1/2} \Sigma = Q_1 \Sigma^{1/2} \end{aligned}$$

- **Cross-covariance linked with Q_1 !**
- Choosing suitable cross-covariance allows the selection of a “good” Q_1 (and hence W)

2. **Cross-correlation** $\Psi = P_{z,x}$ between z and x

$$\begin{aligned} \Psi = \text{Cor}(z, x) &= \Phi V^{-1/2} = W \Sigma V^{-1/2} \\ &= Q_2 P^{-1/2} V^{-1/2} \Sigma V^{-1/2} = Q_2 P^{1/2} \end{aligned}$$

- **Cross-correlation linked with Q_2 !**
- Choosing suitable cross-correlation allows the selection of a “good” Q_2 (and hence W)

Properties:

$$\mathbf{\Psi} = (\psi_{ij})$$

$$\mathbf{\Psi}^T \mathbf{\Psi} = \mathbf{P}$$

$$\implies \text{Diag}(\mathbf{\Psi}^T \mathbf{\Psi}) = (1, 1, \dots, 1) \text{ and } \text{Tr}(\mathbf{\Psi}^T \mathbf{\Psi}) = d$$

$$\implies \sum_{i=1}^d \psi_{ij}^2 = 1$$

i.e. the *column* sums of ψ_{ij}^2 are equal to 1.

Interpretation: this is the squared multiple correlation coefficient $R^2 = 1$ between z_1, \dots, z_d and x_j ! (as the z_i are all uncorrelated you can simply sum the squared correlations to get R^2 ; it is equal to 1 because whitening is an invertible transformation)

The R^2 -s going from x_1, \dots, x_d to the z_i also equal 1, but because the x_i are correlated they need to be computed as $\text{Diag}(\mathbf{\Psi} \mathbf{P}^{-1} \mathbf{\Psi}^T) = \text{Diag}(\mathbf{P}_{z,x} \mathbf{P}^{-1} \mathbf{P}_{x,z}) = (1, 1, \dots, 1)$.

2.4 Natural whitening procedures

Now we discuss several strategies (maximise correlation between individual components, maximise compression, etc.) to arrive at optimal whitening transformation.

This leads to the following “natural” whitening transformations:

- **Mahalanobis** whitening (also known as **ZCA** whitening in machine learning)
- **ZCA-cor** whitening
- **Cholesky** whitening
- **PCA** whitening
- **PCA-cor** whitening
- **CCA** whitening (Canonical Correlation Analysis)

In the following $x_c = x - \mu_x$ and $z_c = z - \mu_z$ denote the mean-centered variables.

2.4.1 ZCA Whitening

Aim: remove correlations but otherwise keep z as similar as possible to x (componentwise!)

$$\begin{aligned} z_1 &\leftrightarrow x_1 \\ z_2 &\leftrightarrow x_2 \\ z_3 &\leftrightarrow x_3 \\ &\vdots \end{aligned}$$

Objective function: minimise $E((z_c - x_c)^T(z_c - x_c)) = d - 2\text{Tr}(\mathbf{\Phi}) + \text{Tr}(\mathbf{V})$

Equivalent objective: maximise $\text{Tr}(\Phi) = \text{Tr}(Q_1 \Sigma^{1/2})$ to find optimal Q_1

Optimal solution:

$$Q_1^{\text{ZCA}} = I$$

$$\implies W^{\text{ZCA}} = \Sigma^{-1/2}$$

- ZCA/Mahalanobis transform is the unique transformation that minimises the expected squared component-wise difference between x and z .
- Use ZCA and Mahalanobis whitening if you want to “just” remove correlations.

2.4.2 ZCA-Cor Whitening

Aim: same as above but remove scale in x first before comparing to z

Objective function: minimise $E \left((z_c - V^{-1/2} x_c)^T (z_c - V^{-1/2} x_c) \right) = 2d - 2\text{Tr}(\Psi)$

Equivalent objective: maximise $\text{Tr}(\Psi) = \text{Tr}(Q_2 P^{1/2})$ to find optimal Q_2

Optimal solution:

$$Q_2^{\text{ZCA-Cor}} = I$$

$$\implies W^{\text{ZCA-Cor}} = P^{-1/2} V^{-1/2}$$

- ZCA-cor whitening is the unique whitening transformation if you aim to maximise correlation between corresponding components in x and z .
- Only if x is standardised to $\text{Var}(x_i) = 1$ then ZCA and ZCA-cor are identical
- ZCA and ZCA-cor: lead to interpretable z

2.4.3 Cholesky Whitening

Aim: find a whitening transformation such that the cross-covariance and cross-correlation have triangular structure. This is useful in some models (such as time course data) to ensure that the future cannot influence the past.

Solution: Cholesky decomposition of $\Sigma^{-1} = LL^T$

L is a lower triangular matrix with positive diagonal elements

$$\implies W^{\text{Chol}} = L^T$$

By construction, W^{Chol} satisfies the whitening constraint since $(W^{\text{Chol}})^T W^{\text{Chol}} = \Sigma^{-1}$.

The corresponding rotation matrices are $Q_1^{\text{Chol}} = L^T \Sigma^{1/2}$ and $Q_2^{\text{Chol}} = L^T V^{1/2} P^{1/2}$

which results in $\Phi^{\text{Chol}} = L^T \Sigma$ and $\Psi^{\text{Chol}} = L^T \Sigma V^{-1/2}$

2.4.4 PCA Whitening

Aim: remove correlations and compress information in \mathbf{x} (each component z_i maximally linked with all variables in \mathbf{x}):

$$\begin{array}{llll} z_1 & \leftarrow & x_1 & z_2 & \leftarrow & x_1 & \dots \\ z_1 & \leftarrow & x_2 & z_2 & \leftarrow & x_2 & \\ \vdots & & & & & & \\ z_1 & \leftarrow & x_d & z_2 & \leftarrow & x_d & \end{array}$$

Objective: maximise $\sum_{j=1}^d \text{Cov}(z_i, x_j)^2 = \sum_{j=1}^d \phi_{ij}^2$ for all i

Equivalent objective: maximise $(\phi_1, \dots, \phi_d)^T = \text{Diag}(\Phi \Phi^T) = \text{Diag}(Q_1 \Sigma Q_1^T)$, with $\phi_1 > \phi_2 > \dots > \phi_d$

Optimal solution:

$$\begin{aligned} Q_1^{\text{PCA}} &= \mathbf{U}^T \\ \implies \mathbf{W}^{\text{PCA}} &= \mathbf{U}^T \Sigma^{-1/2} = \Lambda^{-1/2} \mathbf{U}^T \end{aligned}$$

- Optimum value: $\text{Diag}(Q_1^{\text{PCA}} \Sigma (Q_1^{\text{PCA}})^T) = \text{Diag}(\Lambda) = (\lambda_1, \dots, \lambda_d)^T$
- PCA whitening is the unique whitening transformation that maximises compression with the sum of squared cross-covariances as underlying optimality criterion
- One may use $\frac{\lambda_i}{\sum_{j=1}^d \lambda_j}$ as measure of “proportion of explained variation” for each component in \mathbf{z}
→ scree plots
- This allows the reduction of dimension by discarding low ranking components in \mathbf{z} that do not carry much information

2.4.5 PCA-cor Whitening

Aim: as for PCA whitening but remove scale in \mathbf{x} first (i.e. use squared correlations rather squared covariances to measure compression)

Objective: maximise $\sum_{j=1}^d \text{Cor}(z_i, x_j)^2 = \sum_{j=1}^d \psi_{ij}^2$ for all i

Equivalent objective: maximise $(\psi_1, \dots, \psi_d)^T = \text{Diag}(\Psi \Psi^T) = \text{Diag}(Q_2 P Q_2^T)$ with $\psi_1 > \psi_2 > \dots > \psi_d$.

Optimal solution:

$$Q_2^{\text{PCA-Cor}} = \mathbf{G}^T$$

(with eigendecomposition of $\mathbf{P} = \mathbf{G} \Theta \mathbf{G}^T$)

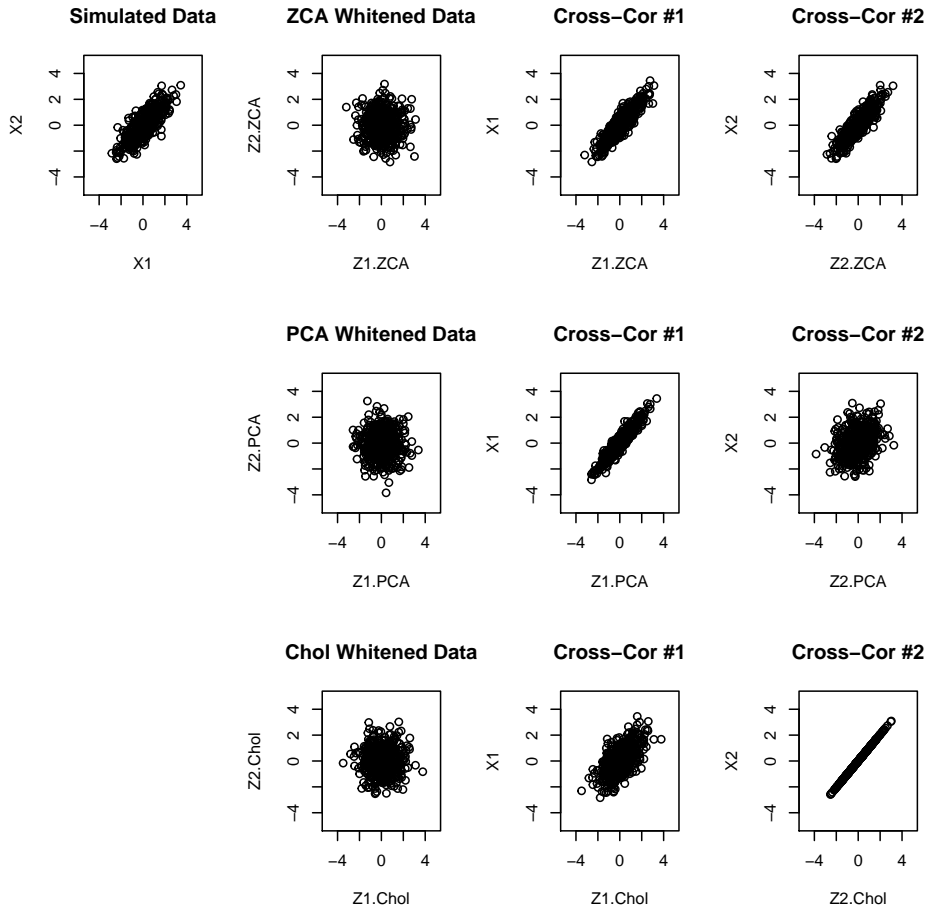
$$\implies \mathbf{W}^{\text{PCA-Cor}} = \Theta^{-1/2} \mathbf{G}^T \mathbf{V}^{-1/2}$$

- Optimum value: $\text{Diag}(Q_2^{\text{PCA-Cor}} \mathbf{P} (Q_2^{\text{PCA-Cor}})^T) = \text{Diag}(\Theta) = (\theta_1, \dots, \theta_d)^T$
- PCA-cor whitening is the unique whitening procedure that maximises compression with the sum of squared cross-correlation as optimality criterion
- If \mathbf{x} is standardised to $\text{Var}(x_i) = 1$, then PCA and PCA-cor are identical.

- $\frac{\theta_i}{\sum_{j=1}^d \theta_j} = \frac{\theta_i}{d}$ indicates relative importance of each whitened variable z_j to predict x
(note that $\sum_{j=1}^d \theta_j = \text{Tr}(\mathbf{P}) = d$)
- Useful for reduction of dimension by discarding low-ranking elements in z

Criterion for relative contributions is linked to multivariate regression R^2 : $\text{Tr}(\mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy})$ which for whitening transformations is d in both directions.

2.4.6 Comparison of ZCA, PCA and Chol whitening



In the above comparison you see ZCA, PCA and Cholesky whitening applied to a simulated bivariate normal data set with correlation $\rho = 0.8$ (column 1). All approaches equally succeed in whitening (column 2) but they differ in the cross-correlations. Columns 3 and 4 shows the cross-correlations between the first two pairs corresponding components for ZCA, PCA and Cholesky whitening. As expected, in ZCA both components show strong correlation, but this is not the case for PCA and Cholesky whitening.

2.4.7 Recap

Method	Type of usage
ZCA, ZCA-cor:	pure decorrelate, maintain similarity to original data set, interpretability
PCA, PCA-cor:	compression, find effective dimension, reduce dimensionality, feature identification
Chol:	time course data

Related models not discussed in this course:

- Factor models: essentially whitening plus an additional error term, factors have rotational freedom just like in whitening
- PLS: similar to PCA but in regression setting (with the choice of latent variables depending on the response)

2.5 Principal Component Analysis (PCA)

- Traditional PCA (invented 1901 by Pearson) is very closely related to **PCA whitening**.
- But PCA itself is **not** a whitening procedure!

PCA transformation:

$$\underbrace{t}_{\text{Principal Components}} = \underbrace{U^T x}_{\text{Orthogonal projection}}$$

$$\implies \text{Var}(t) = \Lambda = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_d \end{pmatrix}$$

- Principal components are **orthogonal** but do *not* have unit variance!
- The variances of the principal components are equal to the eigenvalues of Σ : $\text{Var}(t^{\text{PCA}}) = \Lambda$.
- You arrive at PCA whitening by standardising the PCA components: $z^{\text{PCA}} = \Lambda^{-1/2} t^{\text{PCA}}$
- Same compression / optimality properties as PCA whitening.

2.5.1 Application to data

Written in terms of a data matrix X instead of a random vector x PCA becomes:

$$\underbrace{T}_{\text{Sample version of principal components}} = \underbrace{X}_{\text{Data matrix}} U$$

There are now two ways to obtain U :

1. Estimate the covariance matrix, e.g. by $\hat{\Sigma} = \frac{1}{n} X_c^T X_c$ where X_c is the column-centred data matrix; then apply the eigenvalue decomposition on $\hat{\Sigma}$ to get U .

2. Compute the singular value decomposition of $X_c = VDU^T$. As $\hat{\Sigma} = \frac{1}{n}X_c^T X_c = U(\frac{1}{n}D^2)U^T$ you can just use U from the SVD of X_c and there is no need to compute the covariance.

2.6 PCA correlation loadings and plot

A useful quantity to evaluate the PCA whitened components z^{PCA} and the related principle components t^{PCA} is the cross-correlation matrix Ψ . Specifically, $\Psi = \text{Cor}(z^{\text{PCA}}, x) = \text{Cor}(t^{\text{PCA}}, x) = \Lambda^{1/2}U^T V^{-1/2}$.

We now consider the back-transformation of the (standardised) PCA components to the standardised original components. The inverse PCA transformation is $x = Ut^{\text{PCA}}$ (note that $U^{-1} = U^T$). If one standardises x this equation becomes $V^{-1/2}x = V^{-1/2}Ut^{\text{PCA}}$ and expressed in terms of the standardised z^{PCA} it becomes $V^{-1/2}x = (V^{-1/2}U\Lambda^{1/2})\Lambda^{-1/2}t^{\text{PCA}}$ which is equivalent to $V^{-1/2}x = \Psi^T z^{\text{PCA}}$. Thus the cross-correlation matrix Ψ plays the role of the *correlation loadings*, i.e. the coefficients linking the (standardised) PCA components with the standardised original components.

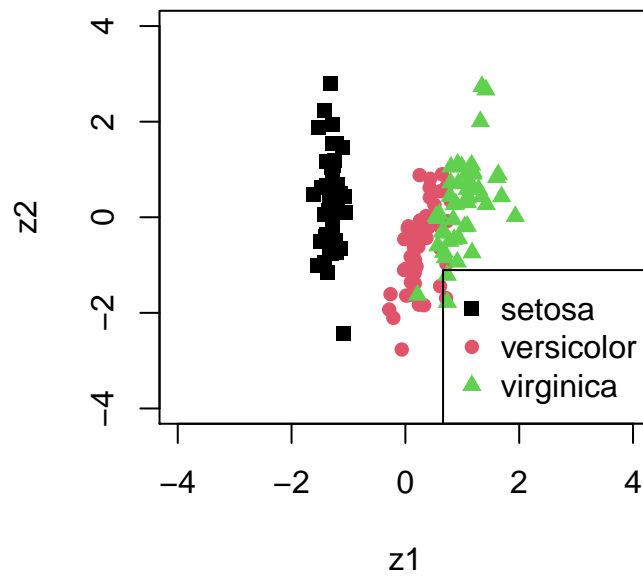
Recall that in a linear regression model with *uncorrelated predictors* the (squared) correlations between each predictor and the response can be used as measure of variable importance. Since PCA (whitened) components are uncorrelated by construction, the correlation loadings Ψ measure the capability of each principal component to predict the original variables.

For visualisation, a correlation loadings plot is often constructed as follows. The loadings Ψ for the first two whitened PCA components z_1, z_2 (or equivalently for t_1, t_2) to all original variables are computed. Then a point for each original variable is drawn in a plane with the two correlation loadings acting as its coordinates. By construction, all points have to lie within a unit circle around the origin. The original variables most strongly influenced by the two latent variables will have strong correlation and thus lie near the outer circle, whereas variables that are not influenced by the two latent variables will lie near the origin. (see next section for an example).

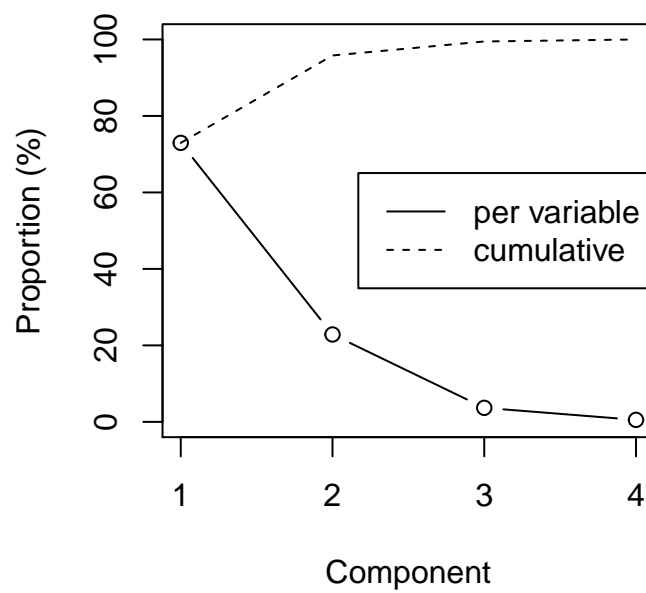
2.6.1 Iris data example

A plot of the the first two components after PCA-whitening is applied reveals the group structure among iris flowers:

PCA Whitening – Iris Data

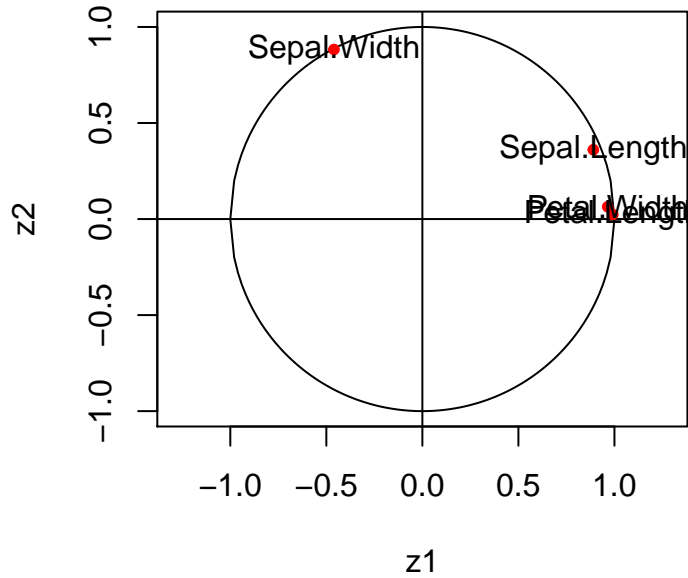


Proportion of Total Variation



Here is the corresponding loadings plot:

Correlation Loadings Iris Data



2.7 CCA whitening (Canonical Correlation Analysis)

So far, we have looked only into whitening as a **single** vector x . In CCA whitening we consider **two vectors** x and y simultaneously:

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_q \end{pmatrix} \quad \begin{array}{l} \text{Var}(x) = \Sigma_x = V_x^{1/2} P_x V_x^{1/2} \\ \text{Var}(y) = \Sigma_y = V_y^{1/2} P_y V_y^{1/2} \end{array}$$

Dimension p Dimension q

$$\text{Whitening of } x: z_x = W_x x = Q_x P_x^{-1/2} V_x^{-1/2} x$$

$$\text{Whitening of } y: z_y = W_y y = Q_y P_y^{-1/2} V_y^{-1/2} y$$

(note we use the correlation-based form of W)

Cross-correlation between z_y and z_x :

$$\text{Cor}(z_x, z_y) = Q_x K Q_y^T$$

with $K = P_x^{-1/2} P_{xy} P_y^{-1/2}$.

Idea: we can choose a suitable Q_x and Q_y rotation matrix by putting constraints on the cross-correlation.

CCA: we aim for a *diagonal* $\text{Cor}(z_x, z_y)$ so that each component in z_x only influences one (the corresponding) component in z_y .

$$z_x = \begin{pmatrix} z_1^x \\ z_2^x \\ \vdots \\ z_p^x \end{pmatrix} \quad z_y = \begin{pmatrix} z_1^y \\ z_2^y \\ \vdots \\ z_q^y \end{pmatrix} \quad \begin{array}{l} \textbf{Motivation} \\ \text{pairs of "modules" represented by components of } z_x \text{ and } z_y \\ \text{influencing each other (and not anyone else).} \end{array} \quad (2.1)$$

$$\text{Cor}(z_x, z_y) = \begin{pmatrix} d_1 & \dots & 0 \\ \vdots & \vdots & \\ 0 & \dots & d_m \end{pmatrix}$$

where d_i are the *canonical correlations* and $m = \min(p, q)$.

2.7.1 How to make cross-correlation matrix $\text{Cor}(z_x, z_y)$ diagonal?

- Use Singular Value Decomposition (SVD) of matrix K :

$$K = (Q_x^{\text{CCA}})^T D Q_x^{\text{CCA}}$$

where D is the diagonal matrix containing the singular values of K

- This yields rotation matrices Q_x^{CCA} and Q_y^{CCA} and thus the desired whitened matrices W_x^{CCA} and W_y^{CCA}
- As a result $\text{Cor}(z_x, z_y) = D$ i.e. singular values of K are identical to canonical correlations d_i !

→ Q_x^{CCA} and Q_y^{CCA} are uniquely determined by the diagonality constraint (and different to the other previously discussed whitening methods)

2.7.2 Related methods

- O2PLS: similar to CCA but using orthogonal projections (thus in O2PLS the latent variables underlying x and y are not orthogonal)

Chapter 3

Clustering / unsupervised Learning

3.1 Overview of clustering

3.1.1 General aim

Find structure and gain insights in data x_1, \dots, x_n collected on n objects by **categorising the objects into groups** based on the d features (= the d components in x_i) obtained for each object.

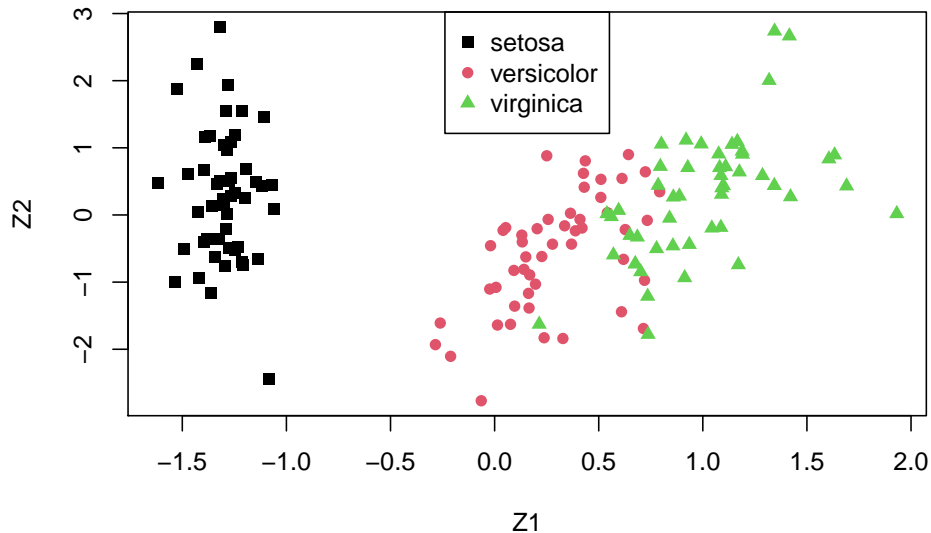
In machine learning / statistical learning this process is called *unsupervised learning* or *clustering*, and there are many algorithms and procedures to automatise and quantify this process.

Unsupervised learning is a very hard problem!

Note that **unsupervised learning** the class labels are a priori unknown, and are learned in the process of clustering. In contrast, in **supervised learning** the class labels are known, at least for the training data set.

Recalle the Isis flower data (from Worksheet~4): PC1 vs. PC2 shows clear visual grouping into 2-3 clusters:

PCA Whitening – Iris Data



Thus, by clustering this data we aim to identify this structure that is visualised here by the *given* class labels.

Two extremes in clustering:

- 1) put all objects into a single cluster (low complexity model)
- 2) put each object into its own cluster (high complexity model)

In both instances nothing new has been learned! In practise, the aim is to find a compromise, i.e. a model that captures the structure in the data with appropriate complexity (not too low and not too complex).

Questions / Problems:

- how do we define clusters?
- how do we learn / infer clusters?
- how many clusters? (this is surprisingly difficult!)
- which features define / separate each cluster?
- uncertainty of clusters?

⇒ Clustering / partitioning / structure discovery is not easy!

⇒ We cannot expect perfect answers or a single “true” clustering (this is related to the problem of model selections, many different clusterings may fit the data equally well)

⇒ Ideally we would wish to get information about the uncertainty of the clustering solution (e.g. by considering sets of possible clusters)

3.1.2 Why is clustering difficult?

Partitioning problem (combinatorics): How many partitions of n objects (say flowers) into K groups (say species) exists?

Answer:

$$S(n, K) = \left\{ \begin{matrix} n \\ K \end{matrix} \right\}$$

this is the “Sterling number of the second type”.

For large n :

$$S(n, K) \approx \frac{K^n}{K!}$$

Example:

n	K	Number of possible partitions
15	3	≈ 2.4 million (10^6)
20	4	≈ 2.4 billion (10^9)
\vdots		
100	5	$\approx 6.6 \times 10^{76}$

These are enormously big numbers even for relatively small problems!

3.1.3 Common types of clustering methods

There are very many different clustering algorithms!

We consider the following two broad types of methods:

- 1) **Algorithmic clustering methods** (these are not explicitly based on a probabilistic model)

- K -means
- PAM
- hierarchical clustering (distance or similarity-based, divide and agglomerative)

pros: fast, effective algorithms to find at least some grouping

cons: no probabilistic interpretation, blackbox methods

- 2) **Model-based clustering** (based on a probabilistic model)

- mixture models (e.g. Gaussian mixture models, GMMs, non-hierarchical)
- graphical models (e.g. Bayesian networks, Gaussian graphical models GGM, trees and networks)

pros: full probabilistic model with all corresponding advantages

cons: computationally very expensive, sometimes impossible to compute exactly.

3.2 Hierarchical clustering

3.2.1 Tree-like structures

Often, categorisations of objects are naturally nested, i.e. there sub-categories of categories etc. These can be represented by **tree-like hierarchical structures**.

In many branches of science hierarchical clusterings are widely employed, for example in evolutionary biology: see e.g.

- Tree of Life with linking the three natural kingdoms
- phylogenetic trees among species (e.g. vertebrata)
- population genetic trees to describe human evolution
- taxonomic trees for plant species
- etc.

Note that when visualising hierarchical structures typically the corresponding tree is depicted facing downwards, i.e. the root of the tree is shown on the top, and the tips/leaves of the tree are shown at the bottom!

In order to obtain such a hierarchical clustering from data two opposing strategies are commonly used:

- 1) **divisive or recursive partitioning algorithms**
 - grow the tree from the root downwards
 - first determine the main two clusters, then recursively refine the clusters further
- 2) **agglomerative algorithms**
 - grow the tree from the leafs upwards
 - successively form partitions by first joining individual object together, then recursively join groups of items together, until all is merged.

For example, K -means can be turned into a divisive hierarchical clustering algorithm by recursively applying the algorithm with $K = 2$.

In the following we discuss a number of popular hierarchical agglomerative clustering algorithms that are based on the pairwise distances / similarities (a $n \times n$ matrix) among all data points.

3.2.2 Agglomerative hierarchical clustering algorithms

A general algorithm for agglomerative construction of a hierarchical clustering works as follows:

Initialisation:

Compute a dissimilarity / distance matrix between all pairs of objects where “objects” are single data points at this stage but later are also be sets of data points.

Iterative procedure:

- 1) identify the pair of objects with the smallest distance. These two objects are then merged together into a common set. Create an internal node in the tree to describe this coalescent event.

- 2) update the distance matrix by computing the distances between the new set and all other objects. If the new set contains all data points the procedure terminates.

For actual implementation of this algorithm two key ingredients are needed:

- 1) a distance measure $d(\mathbf{a}, \mathbf{b})$ between two data points \mathbf{a} and \mathbf{b} .

This is typically on of the following.

- Euclidean distance $d(\mathbf{a}, \mathbf{b}) = \sqrt{(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})}$
- Manhattan distance $d(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^d |a_i - b_i|$
- maximum norm $d(\mathbf{a}, \mathbf{b}) = \max_{i \in \{1, \dots, d\}} |a_i - b_i|$
- etc

In the end, making the correct choice of distance will require subject knowledge about the data!

- 2) a distance measure between two sets of objects $A = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{n_A}\}$ and $B = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{n_B}\}$ of size n_A and n_B , respectively. The centroids of the two sets is given by $\boldsymbol{\mu}_A = \frac{1}{n_A} \sum_{\mathbf{a}_i \in A} \mathbf{a}_i$ and $\boldsymbol{\mu}_B = \frac{1}{n_B} \sum_{\mathbf{b}_i \in B} \mathbf{b}_i$.

To determine the distance $d(A, B)$ between these two sets the following measures are often employed:

- **complete linkage** (max. distance): $d(A, B) = \max_{\mathbf{a}_i \in A, \mathbf{b}_i \in B} d(\mathbf{a}_i, \mathbf{b}_i)$
- **single linkage** (min. distance): $d(A, B) = \min_{\mathbf{a}_i \in A, \mathbf{b}_i \in B} d(\mathbf{a}_i, \mathbf{b}_i)$
- **average linkage** (avg. distance): $d(A, B) = \frac{1}{n_A n_B} \sum_{\mathbf{a}_i \in A} \sum_{\mathbf{b}_i \in B} d(\mathbf{a}_i, \mathbf{b}_i)$

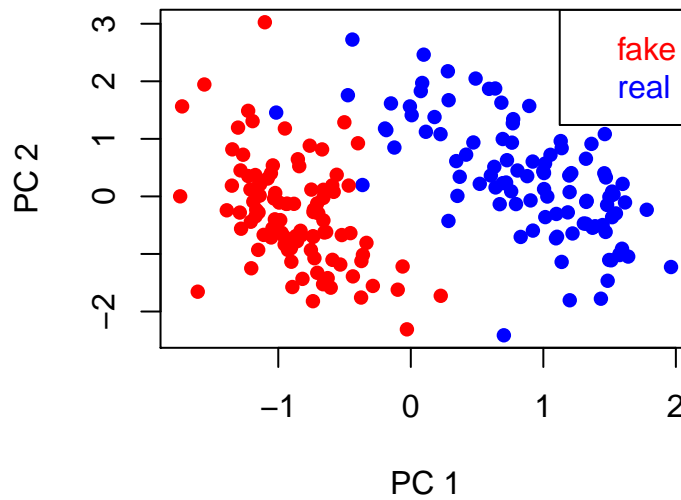
Another agglomerative hierarchical procedure is **Ward's minimum variance approach**. In this approach in each iteration the two sets A and B are merged that lead to the *smallest increase in total within-group sum of squares* (cf. *K-means*). Normally, the within-group A sum of squares $w_A = \sum_{\mathbf{a}_i \in A} (\mathbf{a}_i - \boldsymbol{\mu}_A)^T (\mathbf{a}_i - \boldsymbol{\mu}_A)$ is computed on the basis of actual observations \mathbf{a}_i relative to their mean $\boldsymbol{\mu}_A$. However, it is equally possible to compute it in terms of the squared pairwise differences between the observations using $w_A = \frac{1}{n_A} \sum_{\mathbf{a}_i, \mathbf{a}_j \in A, i < j} (\mathbf{a}_i - \mathbf{a}_j)^T (\mathbf{a}_i - \mathbf{a}_j)$. This is exploited in Ward's clustering method where the distance measure between to sets A and B is $d(A, B) = w_{A \cup B} - w_A - w_B$. Correspondingly, between two data points \mathbf{a} and \mathbf{b} it is the squared Euclidean distance $d(\mathbf{a}, \mathbf{b}) = (\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})$.

3.2.3 Application to Swiss banknote data set

This data set is reports 6 physical measurements on 200 Swiss bank notes. Of the 200 notes 100 are fake and 100 are real. The measurements are: length, left width, right width, bottom margin, top margin, diagonal length of the bank notes.

PCA of this data shows that there are indeed two well defined groups, and that these groups correspond precisely to the real and fake banknotes:

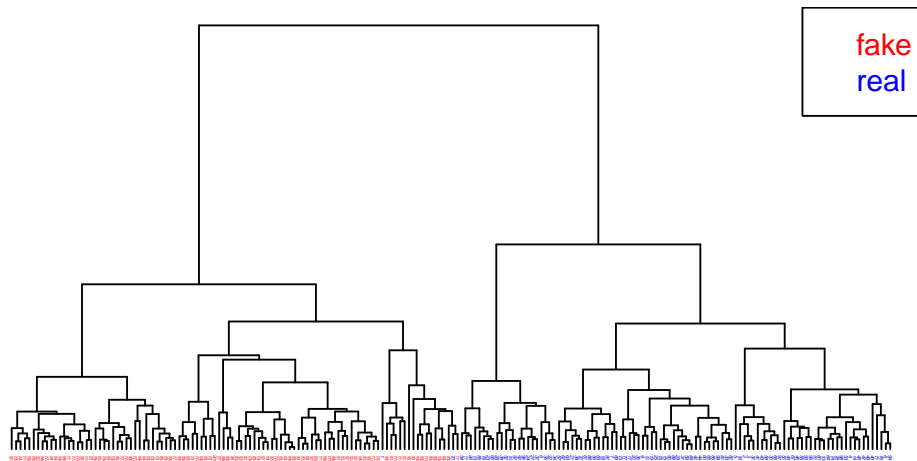
PCA Swiss Banknote Data



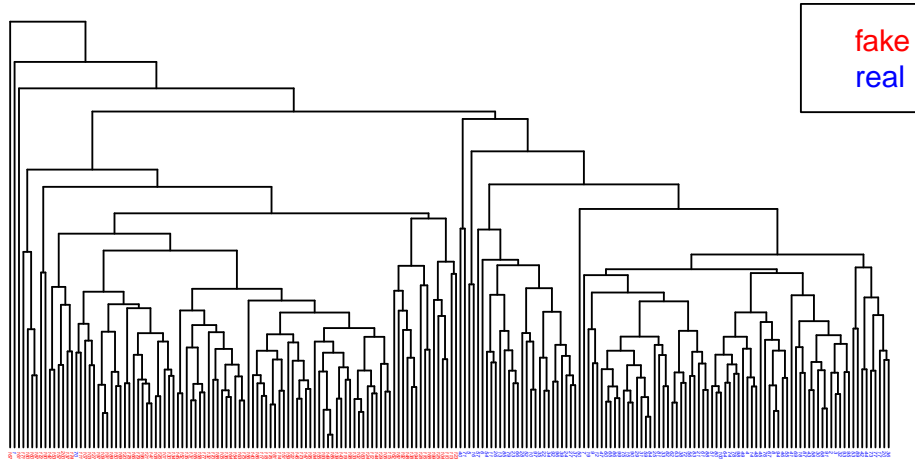
We now compare the clusterings of the above four different hierarchical methods using Euclidean distance:

Ward.D2 (=Ward's method):

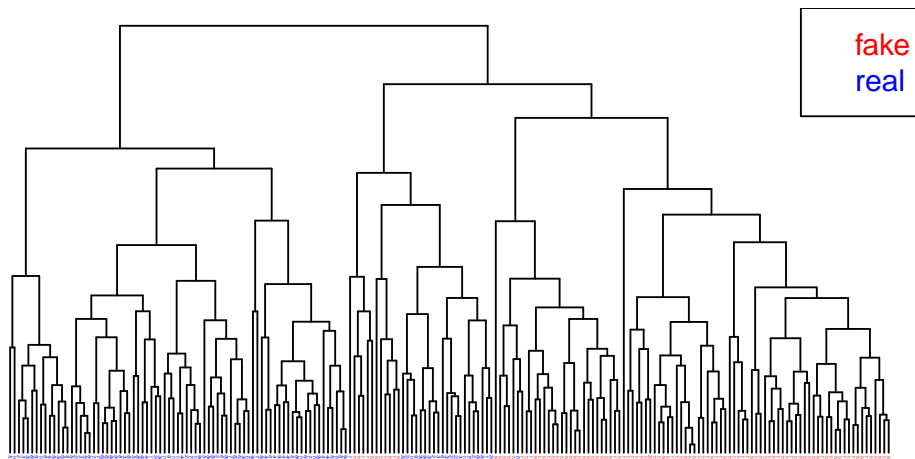
ward.D2 + euclidean



Average linkage:

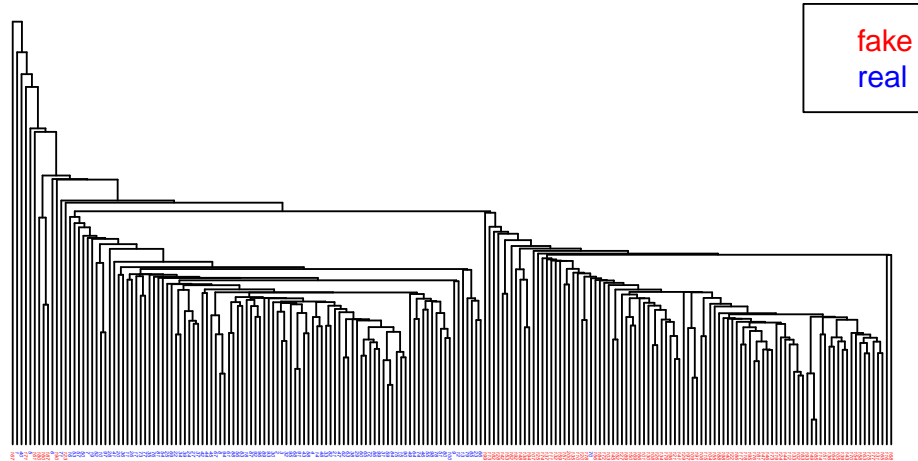
average + euclidean

Complete linkage:

complete + euclidean

Single linkage:

single + euclidean



Result:

- All four trees / hierarchical clusterings are quite different!
- The Ward.D2 method is the only one that finds the correct grouping (except for a single error).

In practical application of hierarchical clustering methods it is essential to evaluate the stability and uncertainty of the obtained groupings. This is often done as follows:

- “bootstrap” (i.e. resampling the original data) is used to generate new data sets (say 200) similar to the original one, and to construct a hierarchical clustering for each of these data sets.
- a consensus tree is computed from the 200 bootstrap trees. This reduces the variability of the estimated tree and also provides bootstrap values measuring the stability of each implied cluster in the tree.

Disadvantage: bootstrapping trees is computationally very expensive!

3.3 K-means clustering

3.3.1 General aims

- Partition the data into K groups, with K given in advance
- The groups are non-overlapping, so each of the n data points / objects x_i is assigned to exactly one of the K groups
- maximise the homogeneity within each group (i.e. each group should contain similar objects)
- maximise the heterogeneity among the different groups (i.e. each group should differ from the other groups)

3.3.2 Algorithm

For each group $k \in \{1, \dots, K\}$ we assume a group mean μ_k . After running K-means we will get estimates of $\hat{\mu}_k$ of the group means, as well as an allocation of each data point to one of the classes.

Initialisation:

At the start of the algorithm the n observations x_1, \dots, x_n are randomly allocated to one of the K groups. The resulting assignment is given by the function $C(x_i) \in \{1, \dots, K\}$. With $G_k = \{i | C(x_i) = k\}$ we denote the set of indices of the data points in cluster k , and with $n_k = |G_k|$ the number of samples in cluster k .

Iterative refinement:

- 1) estimate the group means by

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i \in G_k} x_i$$

- 2) update group assignment: each data point x_i is (re)assigned to the group k with the nearest $\hat{\mu}_k$ (in terms of the Euclidean norm). Specifically, the assignment $C(x_i)$ is updated to

$$\begin{aligned} C(x_i) &= \arg \min_k |x_i - \hat{\mu}_k|_2 \\ &= \arg \min_k (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) \end{aligned}$$

Steps 1 and 2 are repeated until the algorithm converges (or an upper limit of repeats is reached).

3.3.3 Properties

Despite its simplicity K-means is a surprisingly effective clustering algorithms. The final clustering depends on the initialisation so it is often useful to run K-means several times with different starting allocations of the data points.

As a result of the way the clusters are assigned in K-means this leads to cluster boundaries that form a Voronoi tessellation (cf. https://en.wikipedia.org/wiki/Voronoi_diagram) around the cluster means.

Below we will also discuss the connection of K-means with probabilistic clustering using Gaussian mixture models.

3.3.4 Choosing the number of clusters

Once the K-means clustering has been obtained it is insightful to compute:

- a) the total within-group sum of squares SSW (tot.withinss), or total unexplained sum of squares:

$$SSW = \sum_{k=1}^K \sum_{i \in G_k} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k)$$

This quantity decreases with K and is zero for $K = n$. The K -means algorithm tries to minimise this quantity but it will typically only find a local minimum rather than the global one.

- b) the between-group sum of squares SSB (betweeness), or explained sum of squares:

$$SSB = \sum_{k=1}^K n_k (\hat{\mu}_k - \hat{\mu}_0)^T (\hat{\mu}_k - \hat{\mu}_0)$$

where $\hat{\mu}_0 = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{k=1}^K n_k \hat{\mu}_k$ is the global mean of the samples. SSB increases with the number of clusters K until for $K = n$ it becomes equal to

- c) the total sum of squares

$$SST = \sum_{i=1}^n (x_i - \hat{\mu}_0)^T (x_i - \hat{\mu}_0).$$

By construction $SST = SSB + SSW$ for any K (i.e. SST is a constant independent of K).

If divide the sum of squares by the sample size n we get $T = \frac{SST}{n}$ as the *total variation*, $B = \frac{SSB}{n}$ as the *explained variation* and $W = \frac{SSW}{n}$ as the *total unexplained variation*, with $T = B + W$.

For deciding on the optimal number of clusters we can run K -means for various settings of K and then choose the smallest K for which the explained variation B is not significantly worse compared to a model with substantially larger number of clusters (see example below).

3.3.5 K -medoids aka PAM

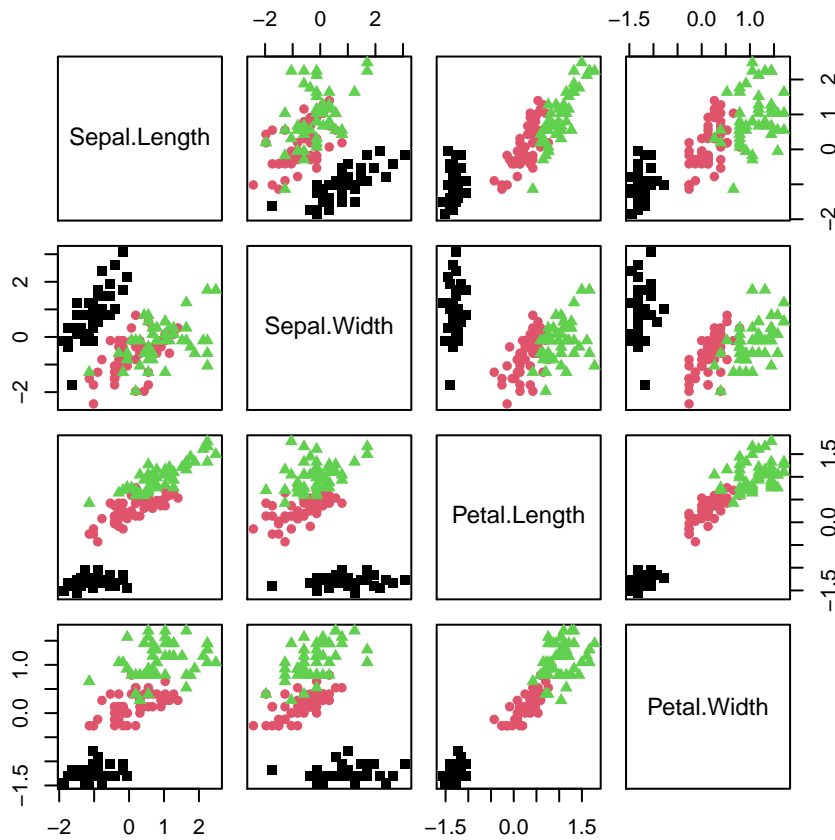
A closely related clustering method is K -medoids or PAM (“Partitioning Around Medoids”).

This works exactly like K -means, only that

- instead of the estimated group means $\hat{\mu}_k$ one member of each group is selected as its representative (the so-called “medoid”)
- instead of squared euclidean distance other dissimilarity measures are also allowed.

3.3.6 Application of K -means to Iris data

Scatter plots of Iris data:



The R output from a *K*-means analysis with true number of clusters specified ($K = 3$) is:

```
## K-means clustering with 3 clusters of sizes 33, 96, 21
##
## Cluster means:
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1   -0.8135055   1.3145538   -1.2825372   -1.2156393
## 2    0.5690971  -0.3705265    0.6888118    0.6609378
## 3   -1.3232208  -0.3718921   -1.1334386   -1.1111395
##
## Clustering vector:
##   [1] 1 3 3 3 1 1 1 1 3 3 1 1 3 3 1 1 1 1 1 1 1 1 1 1 3 1 1 1 3 3 1 1 1 3 3 1
##  [38] 1 3 1 1 3 3 1 1 3 1 3 1 1 2 2 2 2 2 2 3 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2
##  [75] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
## [112] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
## [149] 2 2
##
## Within cluster sum of squares by cluster:
## [1] 17.33362 149.25899 23.15862
## (between_SS / total_SS = 68.2 %)
##
## Available components:
```

```
##
## [1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"
## [6] "betweenss"    "size"         "iter"         "ifault"
```

The corresponding total within-group sum of squares (SSW , `tot.withinss`) and the between-group sum of squares (SSB , `betweenss`) are:

```
kmeans.out$tot.withinss
```

```
## [1] 189.7512
```

```
kmeans.out$betweenss
```

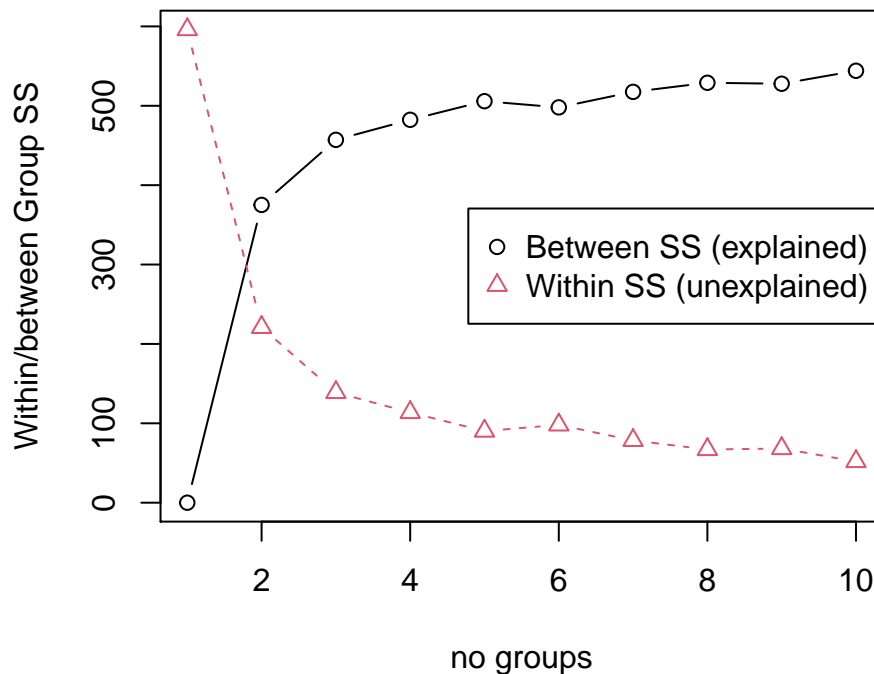
```
## [1] 406.2488
```

By comparing with the known class assignments we can find out about the accuracy of K -means in this example:

```
##
## L.iris      1  2  3
## setosa      33  0 17
## versicolor  0 46  4
## virginica   0 50  0
```

For choosing K we run K -means several times and compute within and between cluster variation in dependence of K :

K-Means Iris Data



Thus, $K = 3$ clusters seem appropriate since the explained variation does not significantly improve (and the unexplained variation does not significantly

decrease) with a further increase of the number of clusters.

3.4 Mixture models

3.4.1 Finite mixture model

- K groups / classes / categories, with the number K specified and finite
- each class $k \in \{1, \dots, K\}$ is modeled by its own distribution F_k with own parameters θ_k .
- density in each class: $f_k(\mathbf{x}) = f(\mathbf{x}|k)$ with $k \in 1, \dots, K$
- mixing weight of each class: $\Pr(k) = \pi_k$ with $\sum_{k=1}^K \pi_k = 1$
- joint density $f(\mathbf{x}, k) = f(\mathbf{x}|k)\Pr(k) = f_k(\mathbf{x})\pi_k$

This results in the mixture density / marginal density

$$f(\mathbf{x}) = \sum_{k=1}^K \pi_k f_k(\mathbf{x})$$

Very often one uses **multivariate normal components** $f_k(\mathbf{x}) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
 \implies **Gaussian mixture model (GMM)**

Mixture models are fundamental not just in clustering but for many other applications (e.g. classification).

Note: don't confuse *mixture model* with *mixed model* (=random effects regression model)

3.4.2 Decomposition of covariance and total variation

Just like in regression you can decompose the variance into an explained and unexplained part.

The conditional means and variances for each class are $E(\mathbf{x}|k) = \boldsymbol{\mu}_k$ and $\text{Var}(\mathbf{x}|k) = \boldsymbol{\Sigma}_k$, and the probability of class k is given by $\Pr(k) = \pi_k$. Using the law of total expectation we can therefore obtain the mean of the mixture density as follows:

$$\begin{aligned} E(\mathbf{x}) &= E(E(\mathbf{x}|k)) \\ &= \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k \\ &= \boldsymbol{\mu}_0 \end{aligned}$$

Similarly, using the law of total variance we compute the marginal variance:

$$\begin{aligned} \underbrace{\text{Var}(\mathbf{x})}_{\text{total}} &= \underbrace{\text{Var}(E(\mathbf{x}|k))}_{\text{explained / between-group}} + \underbrace{E(\text{Var}(\mathbf{x}|k))}_{\text{unexplained / within-group}} \\ \boldsymbol{\Sigma}_0 &= \sum_{k=1}^K \pi_k (\boldsymbol{\mu}_k - \boldsymbol{\mu}_0)(\boldsymbol{\mu}_k - \boldsymbol{\mu}_0)^T + \sum_{k=1}^K \pi_k \boldsymbol{\Sigma}_k \end{aligned}$$

The total variation is given by the trace of the covariance matrix, yielding empirical estimates

$$T = \text{Tr}(\hat{\Sigma}_0) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_0)^T (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_0)$$

$$B = \text{Tr} \left(\sum_{k=1}^K \hat{\pi}_k (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_0) (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_0)^T \right) = \frac{1}{n} \sum_{k=1}^K n_k (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_0)^T (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_0)$$

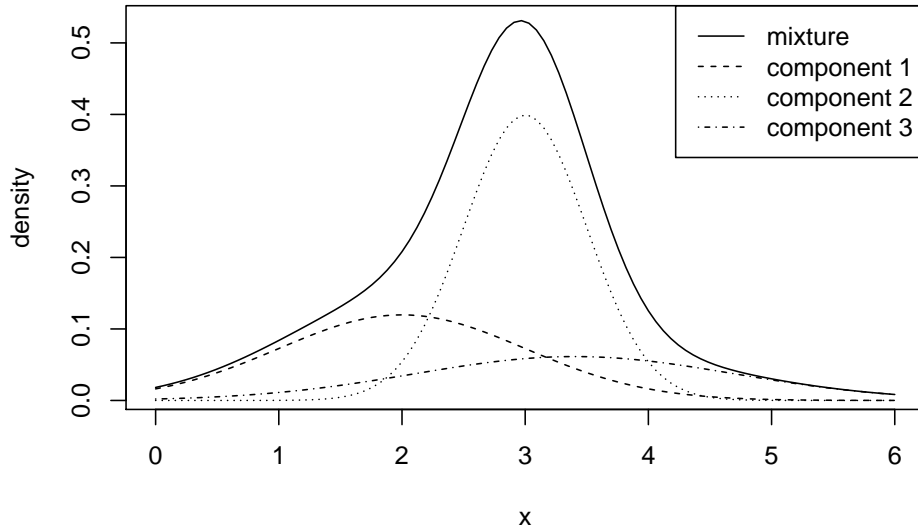
$$W = \text{Tr} \left(\sum_{k=1}^K \hat{\pi}_k \hat{\Sigma}_k \right) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in G_k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)$$

Compare the above with the $T = B + W$ decomposition of total variation in K -means!

3.4.3 Example of mixture of three univariate normal densities:

$$f(x) = 0.3 N(2, 1^2) + 0.5 N(3, 0.5^2) + 0.2 N(3.4, 1.3^2)$$

Gaussian Mixture

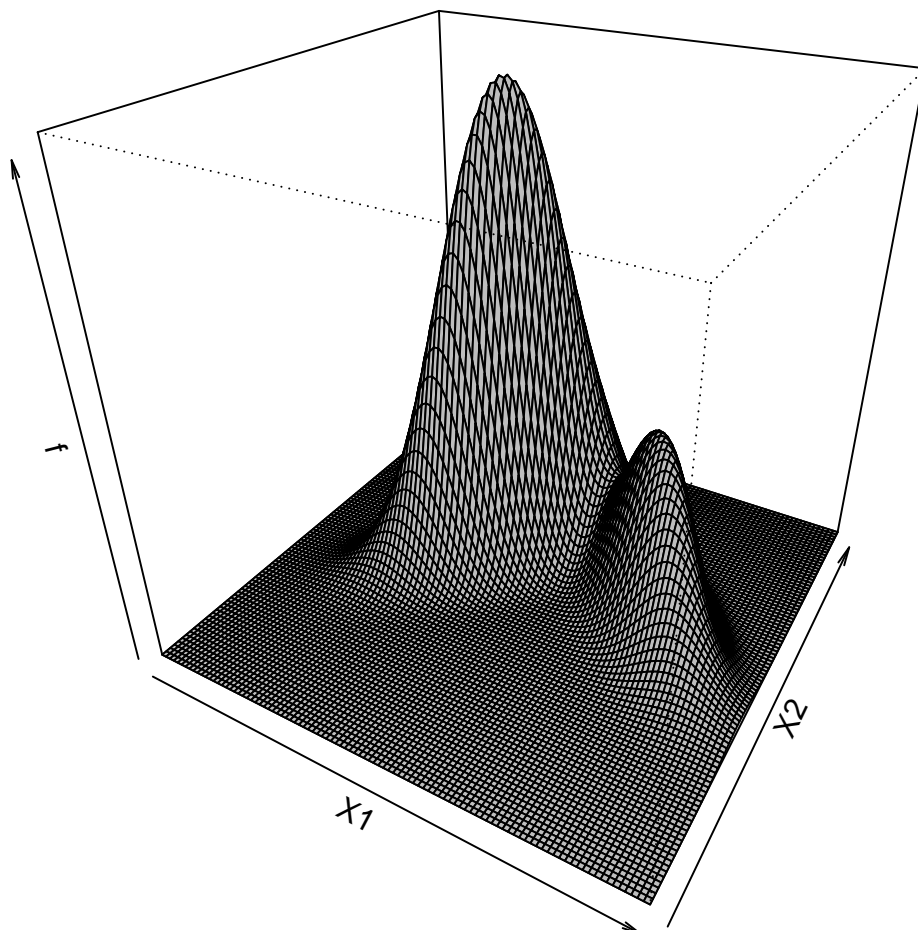


In this case it is clear already by visual inspection that the three subcomponents will not be identifiable.

3.4.4 Example of a mixture of two bivariate normal densities

$$f(\mathbf{x}) = 0.7 N_2 \left(\begin{pmatrix} -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \right) + 0.3 N_2 \left(\begin{pmatrix} 2.5 \\ 0.5 \end{pmatrix}, \begin{pmatrix} 1 & -0.7 \\ -0.7 & 1 \end{pmatrix} \right)$$

Mixture of two bivariate Multinormals



3.4.5 Sampling from a mixture model and latent allocation variable formulation

Assuming we know how to sample from the components $f_k(x)$ of the mixture model it is straightforward to set up a procedure for sampling from the mixture $f(x) = \sum_{k=1}^K \pi_k f_k(x)$.

This is done in a two-step generative process:

1. draw from categorical distribution with parameters $\pi = (\pi_1, \dots, \pi_K)^T$:

$$z \sim \text{Categ}(\pi)$$

the vector $z = (z_1, \dots, z_K)^T$ indicating the group allocation. The group index k is given by $\{k : z_k = 1\}$.

2. Subsequently, sample from the component k selected in step 1:

$$\mathbf{x} \sim F_k$$

This two-stage approach is also called *latent allocation variable formulation* of a mixture model, with \mathbf{z} (or equivalently k) being the latent variable.

The two-step process needs to be repeated for each sample drawn from the mixture (i.e. every time a new latent variable \mathbf{z} is generated).

In probabilistic clustering the aim is to infer the state of \mathbf{z} for all observed samples.

3.4.6 Predicting the group allocation of a given sample

If we know the mixture model and its components we can predict the probability that an observation \mathbf{x} falls in group k using Bayes theorem:

$$z_k = \Pr(k|\mathbf{x}) = \frac{\pi_k f_k(\mathbf{x})}{f(\mathbf{x})}$$

Thus, assuming we can calculate this probability we can **perform probabilistic clustering** by assigning each sample to the class with the largest probability. Unlike in algorithmic clustering, we also get an impression of the uncertainty of the class assignment, since for each sample \mathbf{x} get the vector

$$\mathbf{z} = (z_1, \dots, z_K)^T$$

and thus can see if there are several classes with similar assignment probability. This will be the case, e.g., if \mathbf{x} lies near the boundary between two classes. Note that $\sum_{k=1}^K z_k = 1$.

3.4.7 Variation 1: Infinite mixture model

It is possible to construct mixture models with infinitely many components!

Most commonly known example is Dirichlet process mixture model (DPM):

$$\sum_{k=1}^{\infty} \pi_k f_k(\mathbf{x})$$

with $\sum_{k=1}^{\infty} \pi_k = 1$ and where the weight π_k are taken from a infinitely dimensional Dirichlet distribution (=Dirichlet process).

DPMs are useful for clustering since with them it is not necessary to determine the number of clusters a priori (since it by definition has infinitely many!). Instead, the number of clusters is a by-product of the fit of the model to observed data.

Related: “Chinese restaurant process” - https://en.wikipedia.org/wiki/Chinese_restaurant_process

This describes an algorithm for the allocation process of samples (“persons”) to the groups (“restaurant tables”) in a DPM.

See also “**stick-breaking process**”: https://en.wikipedia.org/wiki/Dirichlet_process#The_stick-breaking_process

3.4.8 Variation 2: Semiparametric mixture model with two classes

A very common model is the following two-component univariate mixture model

$$f(x) = \pi_0 f_0(x) + (1 - \pi_0) f_A(x)$$

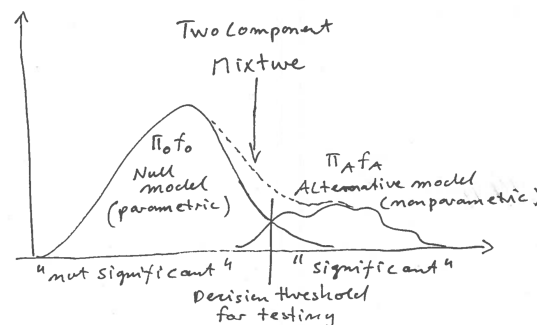
- f_0 : null model, typically parametric such as normal distribution
- f_A : alternative model, typically nonparametric
- π_0 : prior probability of null model

Using Bayes theorem this allows to compute probability that an observation x belongs to the null model:

$$\Pr(\text{Null}|x) = \frac{\pi_0 f_0(x)}{f(x)}$$

This is called the *local false discovery rate*.

The semi-parametric mixture model is the foundation for statistical testing which is based on defining decision thresholds to separate null model (“not significant”) from alternative model (“significant”):



See the lecture notes for Statistical Methods MATH20802 (year 2) for more details.

3.5 Fitting mixture models to data

3.5.1 Direct estimation of mixture model parameters

Given data matrix $X = (x_1, \dots, x_n)^T$ with observations from n samples we would like to fit the mixture model $f(x) = \sum_{k=1}^K \pi_k f_k(x)$ and learn its parameters θ , for example by maximising the corresponding marginal log-likelihood

function with regard to θ :

$$\log L(\theta|X) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k f_k(\mathbf{x}_i) \right)$$

For a Gaussian mixture model the parameters are $\theta = \{\pi, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$.

However, in practise evaluation of this likelihood function may be difficult, in part due to the form of the log-likelihood function (note the sum inside the logarithm), but also due to its singularities and non-identifiability problems.

The above log-likelihood function is also called the *observed data* log-likelihood, or the *incomplete data* log-likelihood, in contrast to the *complete data* log-likelihood described further below.

3.5.2 Estimate mixture model parameters using the EM algorithm

The mixture model may be viewed as an *incomplete* or *missing* data problem: here the missing data are the group allocation $\mathbf{k} = (k_1, \dots, k_n)^T$ belonging to each sample $\mathbf{x}_1, \dots, \mathbf{x}_n$.

If we would know which sample comes from which group the estimation of the parameters θ would indeed be straightforward using the so-called *complete data log-likelihood* based on the joint distribution $f(\mathbf{x}, \mathbf{k}) = f_k(\mathbf{x})\pi_k$

$$\log L(\theta|X, \mathbf{k}) = \sum_{i=1}^n \log (\pi_{k_i} f_{k_i}(\mathbf{x}_i))$$

The idea of the EM algorithm (Dempster et al. 1977) is to exploit the simplicity of the complete data likelihood and to obtain estimates of θ by first finding the probability distribution z_{ik} of the latent variable k_i , and then using this distribution to compute and optimise the corresponding expected complete-data log-likelihood. Specifically, the z_{ik} contain the *probabilities* of each class for each sample i and thus provide a *soft assignment* of classes rather than a 0/1 *hard assignment* (as in the K -means algorithm or in the generative latent variable view of mixture models).

In the EM algorithm we iterate between the

- 1) estimation the probabilistic distribution z_{ik} for the group allocation latent parameters using the current estimate of the parameters θ (obtained in step 2)
- 2) maximisation of the expected complete data log-likelihood to estimate the parameters θ . The expectation is taken with regard to the distribution z_{ik} (obtained in step 1).

Specifically, the EM algorithm applied to model-based clustering proceeds as follows:

- 1) Initialisation: Start with a guess of the parameters $\theta^{(1)}$, then continue with "E" Step, Part A. Alternatively, start with a guess of $z_{ik}^{(1)}$, then continue

with “E” Step, Part B. The initialisation may be derived from some prior information, e.g., from running K -means, or simply be at random.

- 2) **E “expectation” step** — Part A: Use Bayes’ theorem to compute new probabilities of allocation for all the samples \mathbf{x}_i :

$$z_{ik}^{(b+1)} \leftarrow \frac{\pi_k f_k(\mathbf{x}_i)}{f(\mathbf{x}_i)}$$

Note that to obtain $z_{ik}^{(b+1)}$ the current value $\boldsymbol{\theta}^{(b)}$ of the parameters is required.

— Part B: Construct the expected complete data log-likelihood function using the weights $z_{ik}^{(b+1)}$:

$$Q^{(b+1)}(\boldsymbol{\theta}|\mathbf{X}) = \sum_{i=1}^n \sum_{k=1}^k z_{ik}^{(b+1)} \log(\pi_k f_k(\mathbf{x}_i))$$

- 3) **M “maximisation” step** — Maximise the expected complete data log-likelihood to update the mixture model parameters $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}^{(b+1)} \leftarrow \arg \max_{\boldsymbol{\theta}} Q^{(b+1)}(\boldsymbol{\theta}|\mathbf{X})$$

- 4) Repeat with “E” Step until convergence of parameters $\boldsymbol{\theta}^{(b)}$ of the mixture model.

It can be shown that the output $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \boldsymbol{\theta}^{(3)}, \dots$ of the EM algorithm converges to the estimate $\hat{\boldsymbol{\theta}}$ found when maximising the marginal log-likelihood. Since maximisation of the expected complete data log-likelihood is often much easier (and analytically tractable) than maximisation of the observed data log-likelihood function the EM algorithm is the preferred approach in this case.

To avoid singularities in the expected log-likelihood function we may wish to adopt a Bayesian approach (or use regularised/penalised ML) for estimating the parameters in the M-step.

3.5.3 EM algorithm for multivariate normal mixture model

For a GMM the EM algorithm can be written down analytically:

E-step:

$$z_{ik} = \frac{\hat{\pi}_k N(\mathbf{x}_i | \hat{\boldsymbol{\mu}}_k, \hat{\boldsymbol{\Sigma}}_k)}{f(\mathbf{x}_i)}$$

M-step:

$$\begin{aligned} \hat{n}_k &= \sum_{i=1}^n z_{ik} \\ \hat{\pi}_k &= \frac{\hat{n}_k}{n} \end{aligned}$$

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{\hat{n}_k} \sum_{i=1}^n z_{ik} \mathbf{x}_i$$

$$\hat{\boldsymbol{\Sigma}}_k = \frac{1}{\hat{n}_k} \sum_{i=1}^n z_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

Note that the estimators $\hat{\boldsymbol{\mu}}_k$ and $\hat{\boldsymbol{\Sigma}}_k$ are weighted versions of the usual empirical estimators (with weights z_{ik} being the soft assignment of classes resulting from the Bayesian updating).

3.5.4 Connection with K -means clustering method

The K -means algorithm is very closely related to probabilistic clustering with GMMs.

Specifically, it is straightforward to see that *K-means is effectively equivalent to fitting a Gaussian mixture model with the probabilities π_k of all classes identical and with the covariances $\boldsymbol{\Sigma}_k$ all of the form $\sigma^2 \mathbf{I}$* , i.e. all classes have the same diagonal covariance with identical variances. However, note that in K -means the class allocations are hard, whereas in GMMs they are soft. Thus, GMM-based clustering can be viewed as a probabilistic generalisation of K -means clustering!

See also Worksheet~7 where it is shown that the Bayesian update rule assuming equal probability for the classes together with the above covariance leads to the class assignment rule used in K -means.

3.5.5 Choosing the number of classes

Since GMMs operate in a likelihood framework we can use penalised likelihood model selection criteria to choose among different models (i.e. GMMs with different numbers of classes).

The most popular choices are AIC (Akaike Information Criterion) and BIC (Bayesian Information criterion) defined as follows:

$$\text{AIC} = -2 \log L + 2K$$

$$\text{BIC} = -2 \log L + K \log(n)$$

Instead of maximising the log-likelihood we minimise AIC and BIC.

Note that in both criteria more complex models with more parameters (in this case groups) are penalised over simpler models in order to prevent overfitting.

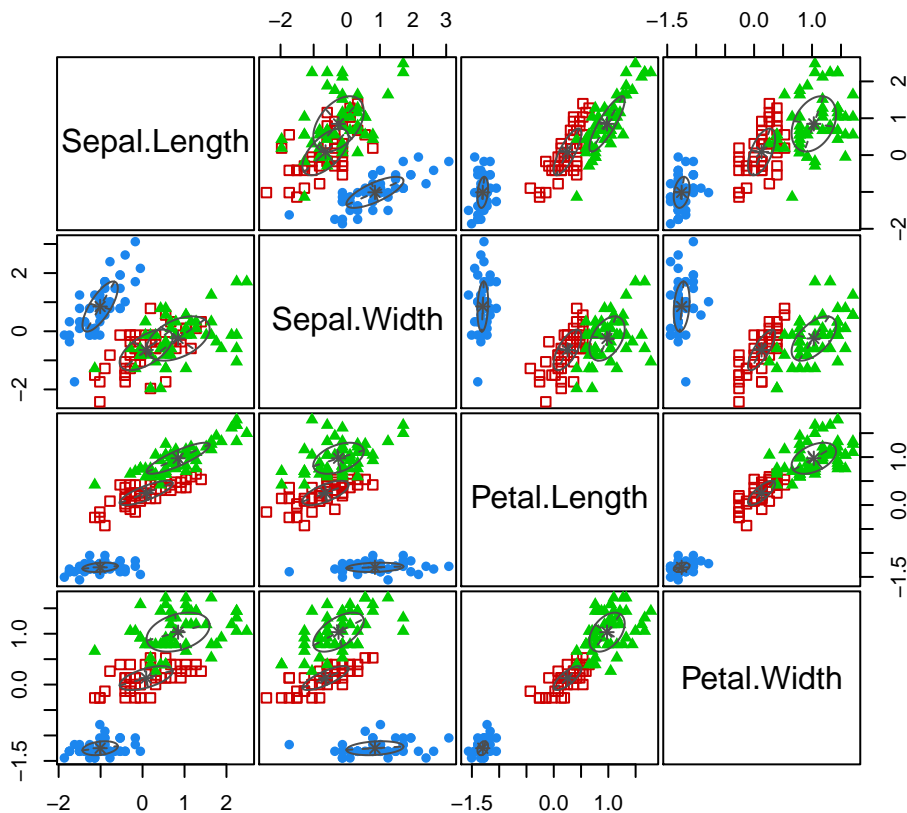
\implies find optimal number of groups K .

Another way of choosing optimal numbers of clusters is by cross-validation (see later chapter on supervised learning).

3.5.6 Application of GMMs to Iris flower data

We now explore the application of GMMs to the Iris flower data set we also investigated with PCA and K-means.

First, we run GMM with 3 clusters:



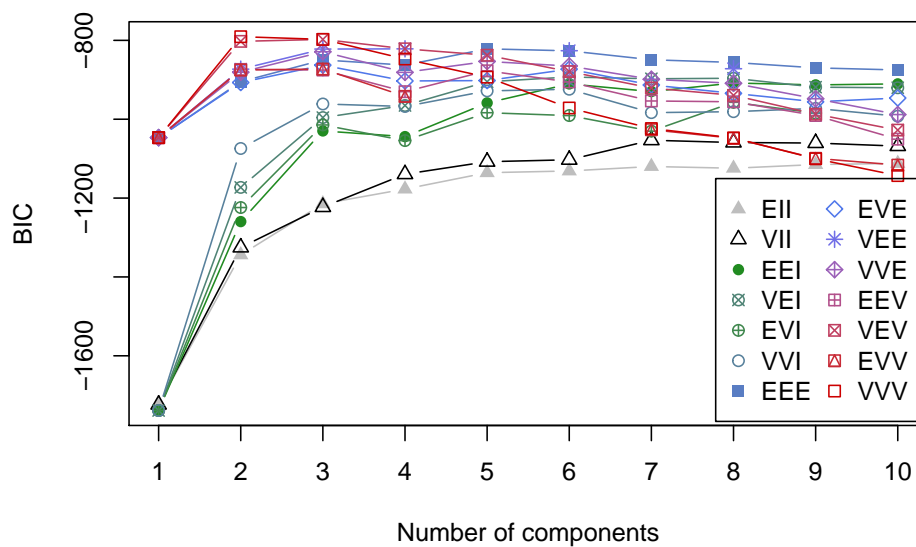
The GMM has a substantially lower misclassification error compared to *K*-means with the same number of clusters:

```
table(gmm3$classification, L.iris)
```

```
##      L.iris
##      setosa versicolor virginica
## 1      50           0           0
## 2       0          45           0
## 3       0           5          50
```

Note that in the R software “mclust” to analyse GMMs the BIC criterion is defined with the opposite sign ($\text{BIC}_{\text{mclust}} = 2 \log L - K \log(n)$), thus we need to find the *maximum* value rather than the smallest value.

If we optimise BIC we find that the model with highest $\text{BIC}_{\text{mclust}}$ is a model with 2 clusters but the model with 3 cluster has nearly as good a BIC:



Chapter 4

Classification / supervised learning

4.1 Introduction

4.1.1 Supervised learning vs. unsupervised learning

Aim in **Unsupervised learning**:

For data x_1, \dots, x_n find classes / labels groups y_1, \dots, y_n attached to each sample x_i .

For example, if x_2 is assigned the label $y = 5$ this means sample 2 belongs to class 5.

If y is discrete unsupervised learning is called *clustering*.

Aim in **Supervised learning**:

We have *training data* available *with* labels: $\{x_1^{train}, y_1^{train}\}, \dots, \{x_n^{train}, y_n^{train}\}$. Each $x_i^{train} = (x_{i1}^{train}, \dots, x_{id}^{train})^T$ contains the observations of d properties (predictor variables) of the sample i .

The training data (observations of predictors variables and response/labels) are used to determine a predictor function $f(x)$. This function is then used to predict the *unknown* labels / class y^{test} of new data x^{test} in a probabilistic fashion, i.e. with probabilities attached to the predicted outcome.

Thus, in contrast to unsupervised learning, supervised learning includes a training step with actual data with known labels.

For y discrete supervised learning is called *classification*.

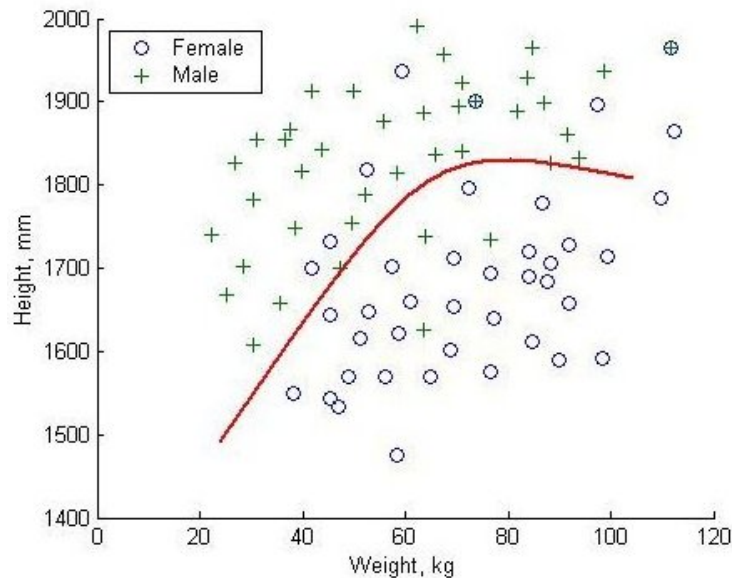
Note the similarity to regression (especially for continuous response y)! In fact, supervised learning *is* (generalised) regression.

4.1.2 Terminology

The function $f(x)$ that predicts the class y is called a *classifier*. There are many types of classifiers, we focus here primarily on probabilistic classifiers (i.e. those that output the predicted class along with a probability). In supervised learning the classifier is learned from the training data.

The challenge is to find a classifier that explains the current training data well *and* that also generalises well to future unseen data. Note that it is relatively easy to find a predictor that explains the training data but especially in high dimensions (i.e. with many predictors) there is often overfitting and then the predictor does not generalise well!

The classifier describes the decision boundary between the classes:



In general, simple decision boundaries are preferred over complex decision boundaries to avoid overfitting!

Some commonly used probabilistic methods for classifications: QDA (quadratic discriminant analysis), LDA (linear discriminant analysis), DDA (diagonal discriminant analysis), Naive Bayes classification, logistic regression, GPs (Gaussian processes).

Common non-probabilistic methods include: SVM (support vector machine), logistic regression, random forest, neural networks.

Depending on how the classifiers are trained there are many variations of the above methods, e.g. Fisher discriminant analysis, regularised LDA, shrinkage discriminant analysis etc.

4.2 Bayesian discriminant rule or Bayes classifier

Same setup as with mixture models:

- K groups with K prespecified
- each group has its own distribution F_k with own parameters θ_k
- the density of each class is $f_k(\mathbf{x}) = f(\mathbf{x}|k)$.
- prior probability of group k is $\Pr(k) = \pi_k$ with $\sum_{k=1}^K \pi_k = 1$
- marginal density is the mixture $f(\mathbf{x}) = \sum_{k=1}^K \pi_k f_k(\mathbf{x})$

The posterior probability of group k is then

$$\Pr(k|\mathbf{x}) = \frac{\pi_k f_k(\mathbf{x})}{f(\mathbf{x})}$$

The *discriminant function* is the logarithm of the posterior probability:

$$d_k(\mathbf{x}) = \log \Pr(k|\mathbf{x}) = \log(\pi_k) + \log(f_k(\mathbf{x})) - \log(f(\mathbf{x}))$$

Since we use d_k to compare the different classes k we can simplify the discriminant function by dropping all constant terms that do not depend on k - in the above $\log(f(\mathbf{x}))$. Hence we get for the Bayes discriminant function

$$d_k(\mathbf{x}) = \log(\pi_k) + \log(f_k(\mathbf{x}))$$

This provides us with the probability of each class given the test data \mathbf{x} . For subsequent “hard” classification we need to use a decision rule, such as selecting the group \hat{k} for that which the group probability / value of discriminant function is maximised:

$$\hat{k} = \arg \max_k d_k(\mathbf{x}).$$

You have already encountered the Bayes classifier in the EM algorithm to predict the state of the latent variables. In a simplified versions it also plays a role in the K-means algorithm.

The Bayes classifier reduces to the likelihood classifier (see example class 3) if one assumes that prior probabilities π_k do not depend on k (and hence are uniform).

4.3 Normal Bayes classifier

4.3.1 Quadratic discriminant analysis (QDA) and Gaussian assumption

Quadratic discriminant analysis (QDA) is a special case of the Bayes classifier when all densities are multivariate normal with $f_k(\mathbf{x}) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

This leads to the discriminant function for QDA:

$$d_k^{QDA}(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) - \frac{1}{2} \log \det(\boldsymbol{\Sigma}_k) + \log(\pi_k)$$

There are a number of noteworthy things here:

- Again terms are dropped that do not depend on k , such as $-\frac{d}{2} \log(2\pi)$.
- Note the appearance of the Mahalanobis distance between x and μ_k in the last term — recall $d^{\text{Mahalanobis}}(x, \mu | \Sigma) = (x - \mu)^T \Sigma^{-1} (x - \mu)$.
- The **QDA discriminant function is quadratic in x** - hence its name!
This implies that the **decision boundaries for QDA classification are quadratic** (i.e. parabolas in two dimensional settings). Thus **QDA is a non-linear classification method!**

For Gaussian models specifically it can be useful to multiply the discriminant function by -2 to get rid of the factor $-\frac{1}{2}$, but note that in that case we then need to look for the minimum of the simplified discriminant function rather than the maximum:

$$d_k^{\text{QDA}(v2)}(x) = (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \det(\Sigma_k) - 2 \log(\pi_k)$$

In the literature you will find both versions of Gaussian discriminant functions so you need to check carefully which convention is used. In the following we will use the first version only.

4.3.2 Linear discriminant analysis (LDA)

LDA is a special case of QDA, with the assumption of common overall covariance across all groups: $\Sigma_k = \Sigma$.

This leads to a simplified discriminant function:

$$d_k^{\text{LDA}}(x) = -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k)$$

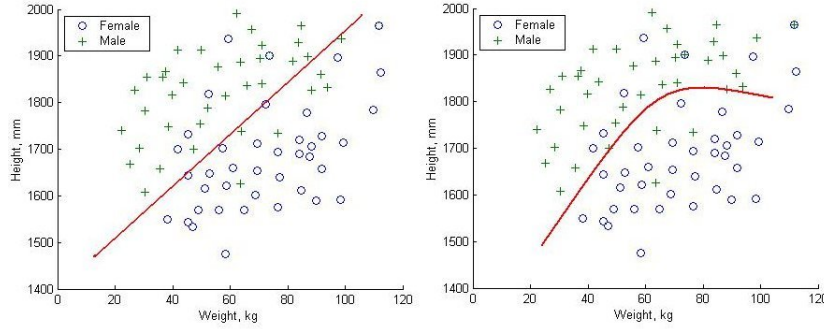
Note that term containing the log-determinant is now gone, and that LDA is essentially now a method that tries to minimize the Mahalanobis distance (while taking also into account the prior class probabilities).

The above function can be further simplified, by noting that the quadratic term $x^T \Sigma^{-1} x$ does not depend on k and hence can be dropped:

$$\begin{aligned} d_k^{\text{LDA}}(x) &= \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\pi_k) \\ &= b^T x + a \end{aligned}$$

Thus, the **LDA discriminant function is linear in x , and hence the resulting decision boundaries are linear** as well (i.e. straight lines in two-dimensional settings). **LDA is a linear classification method.**

Comparison of decision boundary of LDA (left) compared with QDA (right):



Note that logistic regression (cf. GLM module) takes on exactly the above linear form and is indeed closely linked with the LDA classifier.

4.3.3 Diagonal discriminant analysis (DDA)

In DDA we assume the same setting as LDA, but now we simplify even further by assuming a **diagonal covariance** containing only the variances:

$$\Sigma = V = \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{pmatrix}$$

This simplifies the inversion of Σ as

$$\Sigma^{-1} = V^{-1} = \begin{pmatrix} \sigma_1^{-2} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^{-2} \end{pmatrix}$$

and leads to the discriminant function

$$\begin{aligned} d_k^{DDA}(x) &= \mu_k^T V^{-1} x - \frac{1}{2} \mu_k^T V^{-1} \mu_k + \log(\pi_k) \\ &= \sum_{j=1}^d \frac{\mu_{k,j} x_j - \mu_{k,j}^2 / 2}{\sigma_d^2} + \log(\pi_k) \end{aligned}$$

As special case of LDA, the **DDA classifier is a linear classifier**.

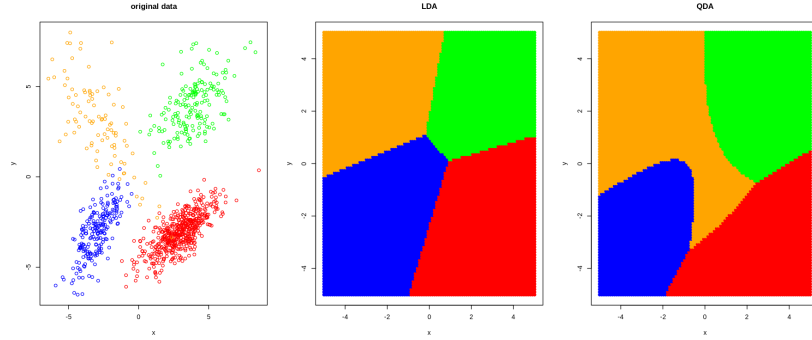
The **Bayes classifier** (using any distribution) **assuming uncorrelated predictors** is also known as the **naïve Bayes classifier**.

Hence, **DDA is a naïve Bayes classifier** assuming underlying Gaussian distributions.

However, don't let you misguide because of the name "naïve": in fact DDA and other "naïve" Bayes classifier are often very effective classifiers, especially in high-dimensional settings!

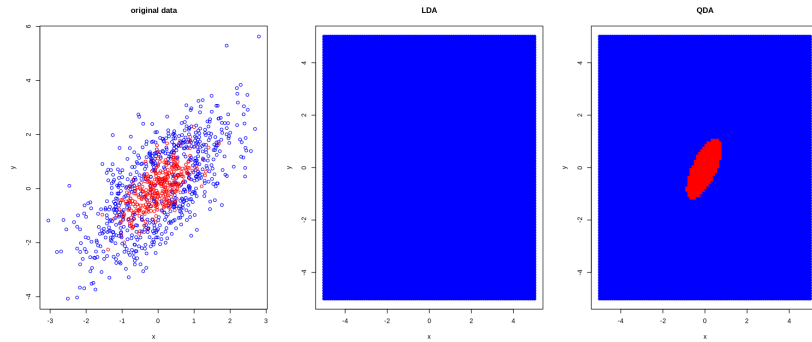
4.3.4 Comparison of decision boundaries: LDA vs. QDA

Non-nested case ($K = 4$)



Note the linear decision boundaries for LDA!

Nested case ($K = 2$):



There is no linear classifier that can separate two nested classes!

4.4 The training step — learning QDA, LDA and DDA classifiers from data

In order to predict the class for new data using any of the above discriminant functions we need to first learn the underlying parameters from the training data $\mathbf{x}_i^{\text{train}}$ and y_i^{train} :

- For QDA, LDA and DDA we need to learn π_1, \dots, π_K .
- For QDA we additionally require $\Sigma_1, \dots, \Sigma_K$
- For LDA we need Σ
- For DDA we estimate $\sigma_1^2, \dots, \sigma_d^2$.

To obtain the above parameter estimates we use the labels y_i^{train} to sort the samples $\mathbf{x}_i^{\text{train}}$ into the corresponding classes, and then apply the usual estimators. Let $G_k = \{i : y_i^{\text{train}} = k\}$ be the set of all indices of training sample belonging to group k .

4.4. THE TRAINING STEP — LEARNING QDA, LDA AND DDA CLASSIFIERS FROM DATA 75

Then to obtain the ML estimate of the group means $k = 1, \dots, K$ we compute

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i \in g_k} x_i^{\text{train}}$$

Note this differs (for $K > 1$) from the estimate of the global mean μ_0 that we get if we were to ignore the group labels (i.e. if we assume there is only a single class):

$$\hat{\mu}_0 = \frac{1}{n} \sum_{i=1}^n x_i^{\text{train}}$$

In order to get the ML estimate of the pooled variance Σ we use

$$\hat{\Sigma}^{ML} = \frac{1}{n} \sum_{k=1}^K \sum_{i \in g_k} (x_i^{\text{train}} - \hat{\mu}_k)(x_i^{\text{train}} - \hat{\mu}_k)^T$$

Note that the pooled variance Σ (with $K > 1$) differs (substantially!) from the global variance Σ_0 that results from ignoring class labels (or in the single class case):

$$\hat{\Sigma}_0^{ML} = \frac{1}{n} \sum_{i=1}^n (x_i^{\text{train}} - \hat{\mu}_0)(x_i^{\text{train}} - \hat{\mu}_0)^T$$

You will recognise the above from the variance decomposition in mixture models, with Σ_0 being the total variance and the pooled Σ the unexplained/within group variance.

Overall, the total number of parameters to be estimated when learning the discriminant functions from training data is as follows:

- QDA: $K + Kd + K \frac{d(d-1)}{2}$
- LDA: $K + d + \frac{d(d-1)}{2}$
- DDA: $K + d$

We also need to make sure that the estimated covariance matrices are all positive definite (which for DDA is automatically guaranteed if all variances are positive).

If d (and K) is small and the number of available samples n is large then we can use maximum likelihood as sketched above to estimate the parameters.

However, if d is large compared to the sample size then the numbers of parameters to estimate grows very quickly. Especially QDA but also LDA is quite data hungry and ML estimation becomes an ill-posed problem. We thus need to use a regularised estimator for the covariance(s) such as penalised ML, Bayes, shrinkage estimator, cf. Section 1.5 and the Statistical Methods module.

Also, to reduce the number of parameters it is advised to use either LDA or DDA in rather than QDA. Often this has a beneficial effect because a simpler model will generalise better and avoid overfitting.

In the application to high-dimensional data we will employ in the computer labs a regularised version of LDA and DDA using the Stein-type shrinkage estimator of the covariance discussed in Section 1.5. Both are implemented in the R package “sda”.

4.5 Goodness of fit and variable selection

As in linear regression (cf. “Statistical Methods” module) we are interested in finding out whether the fitted mixture model is an appropriate model, and which particular predictor(s) x_j from $\mathbf{x} = (x_1, \dots, x_d)^T$ are responsible prediction the outcome, i.e. for categorizing a sample into group k .

In order to study these problem it is helpful to rewrite the discriminant function to highlight the influence (or importance) of each predictor.

We focus on linear methods (LDA and DDA) and first look at the simple case $K = 2$ and then generalise to more than two groups.

4.5.1 LDA with $K = 2$ classes

For two classes using the LDA discriminant rule will choose group $k = 1$ if $d_1^{LDA}(\mathbf{x}) > d_2^{LDA}(\mathbf{x})$, or equivalently, if

$$\Delta_{12}^{LDA} = d_1^{LDA}(\mathbf{x}) - d_2^{LDA}(\mathbf{x}) > 0$$

Since $d_k(\mathbf{x})$ is the log-posterior (plus/minus identical constants) Δ_{12}^{LDA} is in fact the **log-posterior odds of class 1 versus class 2** (see Statistical Methods, Bayesian inference).

The difference Δ_{12}^{LDA} is

$$\underbrace{\Delta_{12}^{LDA}}_{\text{log posterior odds}} = \underbrace{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)}_{\text{log Bayes factor log } B_{12}} + \underbrace{\log \left(\frac{\pi_1}{\pi_2} \right)}_{\text{log prior odds}}$$

Note that since we only consider simple non-composite models here the log-Bayes factor is identical with the log-likelihood ratio!

The log Bayes factor $\log B_{12}$ is known as the *weight of evidence* in favour of F_1 given \mathbf{x} . The *expected weight of evidence* assuming \mathbf{x} is indeed from F_1 is the Kullback-Leibler discrimination information in favour of group 1, i.e. the KL divergence of from distribution F_2 to F_1 :

$$E_{F_1}(\log B_{12}) = KL(F_1 || F_2) = \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = \frac{1}{2}\Omega^2$$

This yields, apart of a scale factor, a population version of the Hotelling T^2 statistic defined as

$$T^2 = c^2(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2)^T \hat{\boldsymbol{\Sigma}}^{-1}(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2)$$

where $c = (\frac{1}{n_1} + \frac{1}{n_2})^{-1/2} = \sqrt{n\pi_1\pi_2}$ is a sample size dependent factor (for $SD(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_2)$). T^2 is a measure of fit of the underlying two-component mixture.

Using the whitening transformation with $\mathbf{z} = \mathbf{W}\mathbf{x}$ and $\mathbf{W}^T\mathbf{W} = \boldsymbol{\Sigma}^{-1}$ we can rewrite the log Bayes factor as

$$\begin{aligned} \log B_{12} &= \left((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{W}^T \right) \left(\mathbf{W} \left(\mathbf{x} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right) \right) \\ &= \boldsymbol{\omega}^T \boldsymbol{\delta}(\mathbf{x}) \end{aligned}$$

i.e. as the product of two vectors:

- $\delta(x)$ is the whitened x (centered around average means) and
- $\omega = (\omega_1, \dots, \omega_d)^T = W(\mu_1 - \mu_2)$ gives the weight of each whitened component $\delta(x)$ in the log Bayes factor.

A large positive or negative value of ω_j indicates that the corresponding whitened predictor is relevant for choosing a class, whereas small values of ω_j close to zero indicate that the corresponding ZCA whitened predictor is unimportant. Furthermore, $\omega^T \omega = \sum_{j=1}^d \omega_j^2 = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) = \Omega^2$, i.e. the squared ω_j^2 provide a component-wise decomposition of the overall fit Ω^2 .

Choosing ZCA-cor as whitening transformation with $W = P^{-1/2} V^{-1/2}$ we get

$$\omega^{ZCA-cor} = P^{-1/2} V^{-1/2} (\mu_1 - \mu_2)$$

A better understanding of $\omega^{ZCA-cor}$ is provided by comparing with the two-sample t -statistic

$$\hat{\tau} = c \hat{V}^{-1/2} (\hat{\mu}_1 - \hat{\mu}_2)$$

With τ the population version of $\hat{\tau}$ we can define

$$\tau^{adj} = P^{-1/2} \tau = c \omega^{ZCA-cor}$$

as correlation-adjusted t -scores (cat scores). With $(\hat{\tau}^{adj})^T \hat{\tau}^{adj} = T^2$ we can see that the cat scores offer a component-wise decomposition of Hotelling's T^2 .

Note the choice of ZCA whitening is to ensure that the whitened components are interpretable and stay maximally correlated to the original variables. However, you may also choose, e.g. PCA whitening, in which case the $\omega^T \omega$ provide the variable importance for the PCA whitened variables.

For DDA, which assumes that correlations among predictors vanish, i.e. $P = I_d$, we get

$$\Delta_{12}^{DDA} = \underbrace{\left((\mu_1 - \mu_2)^T V^{-1/2} \right)}_{c^{-1} \tau^T} \underbrace{\left(V^{-1/2} \left(x - \frac{\mu_1 + \mu_2}{2} \right) \right)}_{\text{centered standardised predictor}} + \log \left(\frac{\pi_1}{\pi_2} \right)$$

Similarly as above, the t -score τ determines the impact of the standardised predictor in Δ^{DDA} .

Consequently, in DDA we can rank predictors by the squared t -score. Recall that in standard linear regression with uncorrelated predictors we can find the most important predictors by ranking the squared marginal correlations – ranking by (squared) t -scores in DDA is the exact analogy but for discrete response.

4.5.2 Multiple classes

For more than two classes we need to refer to the so-called **pooled centroids formulation** of DDA and LDA (introduced by Tibshirani 2002).

We define the pooled centroid as $\mu_0 = \sum_{k=1}^K \pi_k \mu_k$ — this is the centroid if there would be only a single class. The corresponding frequency is $\pi_0 = 1$ and the distribution is called F_0 .

The LDA discriminant function for this “group 0” is

$$d_0^{LDA}(x) = \mu_0^T \Sigma^{-1} x - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0$$

and the log posterior odds for comparison of group k with the pooled group 0 is

$$\begin{aligned} \Delta_k^{LDA} &= d_k^{LDA}(x) - d_0^{LDA}(x) \\ &= \log B_{k0} + \log(\pi_k) \\ &= \omega_k^T \delta_k(x) + \log(\pi_k) \end{aligned}$$

with

$$\omega_k = W(\mu_k - \mu_0)$$

and

$$\delta_k(x) = W(x - \frac{\mu_k + \mu_0}{2})$$

The expected log Bayes factor is

$$E_{F_k}(\log B_{k0}) = KL(F_k || F_0) = \frac{1}{2} (\mu_k - \mu_0)^T \Sigma^{-1} (\mu_k - \mu_0) = \frac{1}{2} \Omega_k^2$$

With scale factor $c_k = (\frac{1}{n_k} - \frac{1}{n})^{-1/2} = \sqrt{n \frac{\pi_k}{1-\pi_k}}$ (for $SD(\hat{\mu}_k - \hat{\mu}_0)$, with the minus sign before $\frac{1}{n}$ due to correlation between $\hat{\mu}_k$ and pooled mean $\hat{\mu}_0$) we get as correlation-adjusted t -score for comparing mean of group k with the pooled mean

$$\tau_k^{adj} = c_k \omega_k^{ZCA-cor}.$$

For the two class case ($K = 2$) we get with $\mu_0 = \pi_1 \mu_1 + \pi_2 \mu_2$ for the mean difference $(\mu_1 - \mu_0) = \pi_2(\mu_1 - \mu_2)$ and with $c_1 = \sqrt{n \frac{\pi_1}{\pi_2}}$ this yields

$$\tau_1^{adj} = \sqrt{n \pi_1 \pi_2} P^{-1/2} V^{-1/2} (\mu_1 - \mu_2),$$

i.e. the exact same score as in the two-class setting.

4.5.3 Choosing a threshold

From the above it is clear that in LDA and DDA the natural score to rank features with regard to importance in the predictor is the (squared) t -score (no correlation) or the squared correlation-adjusted t -score.

In order to determine a suitable threshold one can use any standard technique, such as multiple testing multiple testing or FDR thresholding.

In Computer Lab 4 we will perform feature selection on an example data set using both the t -score and the correlation-adjusted t -score.

This will also show that using feature selection it is often possible to construct compact models with fewer predictors that still generalise and predict well.

For large and high-dimensional models feature selection can also be viewed as a form of regularisation and also dimension reduction. Specifically, if there are many variables/ features that do not contribute to the prediction they can still deteriorate the overall predictive accuracy (sometimes dramatically) so these “noise variables” need to be filtered out in order to be able to construct good models and classifiers.

4.6 Estimating prediction error

4.6.1 Quantifying prediction error

For any prediction model we are interested in the predictive performance. We quantify the performance by comparing the prediction \hat{y} with the true output y (assumed to be known).

For continuous response often the squared loss is used:

$$\text{err}(\hat{y}, y) = (\hat{y} - y)^2$$

For binary outcomes one often employs the 0/1 loss:

$$\text{err}(\hat{y}, y) = \begin{cases} 1, & \text{if } \hat{y} \neq y \\ 0, & \text{otherwise} \end{cases}$$

but we can of course use any other quantity derived from the confusion matrix (containing TP, TN, FP, FN).

The mean prediction error is then the expectation

$$PE = E(\text{err}(\hat{y}, y))$$

and thus the empirical mean prediction error is

$$\widehat{PE} = \frac{1}{m} \sum_{i=1}^m \text{err}(\hat{y}_i, y_i)$$

where m is the sample size of the **test data** (different from the **training data** used to construct the model!!).

Alternatively and more generally, we can also quantify prediction error in the framework of so-called **proper scoring rules**, where the whole probabilistic forecast is taken into account (e.g. the individual probabilities for each class, rather than just the selected most probable class). A commonly used scoring rule is the negative log-probability (“surprise”), and the expected surprise is the cross-entropy (cf. Statistical Methods module). So this leads back to entropy and likelihood.

Once we have an estimate of the prediction error of a model we can use this error to choose among the models (including models with different numbers of features).

4.6.2 Estimation of prediction error without test data

Unfortunately, quite often we do not have any test data available to evaluate a classifier.

In this case we need to rely on a simple algorithmic procedure called **cross-validation**.

Idea:

- split the samples in the training data into a number (say K) parts (“folds”)
- use each of the K folds as test data and the other $K - 1$ folds as training data
- average over the resulting K estimates of prediction error

Note that in each case one part of the data is reserved for testing and not used for training.

We choose K such that the folds are not too small (to allow estimation of prediction error) but also not too large (to make sure that we actually train a good classifier from the remaining data). A typical value for K is 5 or 10.

In Computer labs 4 and 5 we employ cross-validation to estimate prediction accuracy and model selection.

Relevant reading for the technical details: **Section 5.1 Cross-Validation** in [James et al. \(2013\) *An introduction to statistical learning with applications in R*](#). Springer.

Chapter 5

Multivariate dependencies

5.1 Measuring the association between two sets of random variables

5.1.1 Rozeboom vector correlation

In linear regression model the squared multiple correlation (also known as coefficient of determination) between y and x

$$\text{Cor}(y, x)^2 = P_{yx} P_x^{-1} P_{yx}$$

is a standard measure to describe the strength of association between the predictors x and the response y . If there is only a single predictor the squared multiple correlation, or coefficient of determination, reduces to the squared Pearson correlation $\text{Cor}(x, y)^2$.

Now, if we consider two random vectors $x = (x_1, \dots, x_p)^T$ and $y = (y_1, \dots, y_q)^T$, what is a relevant measure to describe the association between x and y that generalises both simple correlation and multiple correlation?

An answer in the form of squared *vector correlation* was given by Rozeboom (1965) defined as follows:

$$\text{Cor}(x, y)^2 \rho_{xy}^2 = 1 - \prod_{i=1}^m (1 - \lambda_i^2)$$

where $\lambda_1, \dots, \lambda_m$ are the canonical correlations, i.e. the singular values of $K = P_x^{-1/2} P_{xy} P_y^{-1/2}$ (cf. Chapter 2).

The Rozeboom vector correlation is the complement of Hotelling's (1936) *vector alienation coefficient* given by

$$a(x, y) = \prod_{i=1}^m (1 - \lambda_i^2)$$

so $\text{Cor}(x, y)^2 = 1 - a(x, y)$.

Equivalent ways to write the squared vector correlation are

$$\begin{aligned}\text{Cor}(x, y)^2 &= 1 - \frac{\det \mathbf{P}_{x,x}}{\det \mathbf{P}_x \det \mathbf{P}_y} \\ &= 1 - \det \left(\mathbf{I}_q - \mathbf{P}_y^{-1} \mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy} \right) \\ &= 1 - \det \left(\mathbf{I}_p - \mathbf{P}_x^{-1} \mathbf{P}_{xy} \mathbf{P}_y^{-1} \mathbf{P}_{yx} \right) .\end{aligned}$$

It is easy to see that for $p = 1$ or $q = 1$ the squared vector correlation reduces to the squared multiple correlation.

Using the Weinstein-Aronszajn determinant identity $\det(\mathbf{I} + \mathbf{AB}) = \det(\mathbf{I} + \mathbf{BA})$ we can see that

$$\begin{aligned}\det \left(\mathbf{I}_q - \mathbf{P}_y^{-1} \mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy} \right) &= \det \left(\mathbf{I}_q - \mathbf{P}_y^{-1/2} \mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy} \mathbf{P}_y^{-1/2} \right) \\ &= \det \left(\mathbf{I}_q - \mathbf{K}^T \mathbf{K} \right)\end{aligned}$$

Thus, the squared vector correlation between x and y written in terms of the matrix \mathbf{K} is

$$\begin{aligned}\rho_{xy}^2 &= 1 - \det \left(\mathbf{I}_q - \mathbf{K}^T \mathbf{K} \right) \\ &= 1 - \det \left(\mathbf{I}_p - \mathbf{K} \mathbf{K}^T \right)\end{aligned}$$

which brings us back to the first definition as the complement of the vector alienation coefficient.

5.1.2 Other common approaches

A common approach to measure association between is the RV coefficient defined as

$$RV(\mathbf{X}, \mathbf{Y}) = \frac{\text{Tr}(\mathbf{\Sigma}_{XY} \mathbf{\Sigma}_{YX})}{\text{Tr}(\mathbf{\Sigma}_X^2) \text{Tr}(\mathbf{\Sigma}_Y^2)}$$

While for $q = p = 1$ the RV coefficient becomes the pairwise squared correlation. However, the RV coefficient does not reduce to the multiple correlation coefficient for $q = 1$ and $p > 1$.

Another way to measure multivariate association is mutual information (MI) which not only covers linear but also non-linear association. See the next Chapter for details. MI applied to multivariate normal distribution is linked to the above Rozeboom vector correlation.

5.2 Graphical models

5.2.1 Purpose

Graphical models combine features from

- graph theory

- probability
- statistical inference

The literature on graphical models is huge, we focus here only on two commonly used models:

- DAGs (directed acyclic graphs), all edges are directed, no directed loops (i.e. no cycles, hence “acyclic”)
- GGM (Gaussian graphical models), all edges are undirected

Graphical models provide probabilistic models for trees and for networks, with random variables represented by nodes in the graphs, and branches representing conditional dependencies. In this regard they generalise both the tree-based clustering approaches as well as the probabilistic non-hierarchical methods (GMMs).

However, the class of graphical models goes much beyond simple unsupervised learning models. It also includes regression, classification, time series models etc. See e.g. the reference book by [Murphy \(2012\)](#).

5.2.2 Basic notions from graph theory

- Mathematically, a graph $G = (V, E)$ consists of a set of vertices or nodes $V = \{v_1, v_2, \dots\}$ and a set of branches or edges $E = \{e_1, e_2, \dots\}$.
- Edges can be undirected or directed.
- Graphs containing only directed edges are directed graphs, and likewise graphs containing only undirected edges are called undirected graphs. Graphs containing both directed and undirected edges are called partially directed graphs.
- A path is a sequence of vertices such that from each of its vertices there is an edge to the next vertex in the sequence.
- A graph is connected when there is a path between every pair of vertices.
- A cycle is a path in a graph that connects a node with itself.
- A connected graph with no cycles is called a tree.
- The degree of a node is the number of edges it connects with. If edges are all directed the degree of a node is the sum of the in-degree and out-degree, which counts the incoming and outgoing edges, respectively.
- External nodes are nodes with degree 1. In a tree-structured graph these are also called leaves.

Some notions are only relevant for graphs with directed edges:

- In a directed graph the parent node(s) of vertex v is the set of nodes $\text{pa}(v)$ directly connected to v via edges directed from the parent node(s) towards v .
- Conversely, v is called a child node of $\text{pa}(v)$. Note that a parent node can have several child nodes, so v may not be the only child of $\text{pa}(v)$.
- In a directed tree graph, each node has only a single parent, except for one particular node that has no parent at all (this node is called the root node).
- A DAG, or directed acyclic graph, is a directed graph with no directed cycles. A (directed) tree is a special version of a DAG.

5.2.3 Probabilistic graphical models

A graphical model uses a graph to describe the relationship between random variables x_1, \dots, x_d . The variables are assumed to have a joint distribution with density/mass function $\Pr(x_1, x_2, \dots, x_d)$. Each random variable is placed in a node of the graph.

The structure of the graph and the type of the edges connecting (or not connecting) any pair of nodes/variables is used to describe the conditional dependencies, and to simplify the joint distribution.

Thus, a graphical model is in essence a visualisation of the joint distribution using structural information from the graph helping to understand the mutual relationship among the variables.

5.2.4 Directed graphical models

In a **directed graphical model** the graph structure is assumed to be a DAG (or a directed tree, which is also a DAG).

Then the joint probability distribution can be factorised into a *product of conditional probabilities* as follows:

$$\Pr(x_1, x_2, \dots, x_d) = \prod_i \Pr(x_i | \text{pa}(x_i))$$

Thus, the overall joint probability distribution is specified by local conditional distributions and the graph structure, with the directions of the edges providing the information about parent-child node relationships.

Probabilistic DAGs are also known as “Bayesian networks”.

Idea: by trying out all possible trees/graphs and fitting them to the data using maximum likelihood (or Bayesian inference) we hope to be able identify the graph structure of the data-generating process.

Challenges

- 1) in the tree/network the internal nodes are usually not known, and thus have to be treated as *latent* variables.

Answer: To impute the states at these nodes we may use the EM algorithm as in GMMs (which in fact can be viewed as graphical models, too!).

- 2) If we treat the internal nodes as unknowns we need to marginalise over the internal nodes, i.e. we need to sum / integrate over all possible set of states of the internal nodes!

Answer: This can be handled very effectively using the **Viterbi algorithm** which is essentially an application of the generalised distributive law. In particular for tree graphs this means that the summations occurs locally at each nodes and propagates recursively accross the tree.

- 3) In order to infer the tree or network structure the space of all trees or networks need to be explored. This is not possible in an exhaustive fashion unless the number of variables in the tree is very small.

Answer: Solution: use heuristic approaches for tree and network search!

- 4) Furthermore, there exist so-called “equivalence classes” of graphical models, i.e. sets of graphical models that share the same joint probability distribution. Thus, all graphical models within the same equivalence class cannot be distinguished from observational data, even with infinite sample size!

Answer: this is a fundamental mathematical problem of identifiability so there is now way around this issue. However, on the positive side, this also implies that the search through all graphical models can be restricted to finding the so-called “essential graph” (e.g. <https://projecteuclid.org/euclid.aos/1031833662>)

Conclusion: using directed graphical models for structure discovery is very time consuming and computationally demanding for anything but small toy data sets.

This also explains why heuristic and non-model based approaches (such as hierarchical clustering) are so popular even though full statistical modelling is in principle possible.

5.2.5 Undirected graphical models

Another class of graphical models are models that contain only undirected edges. These **undirected graphical models** are used to represent the pairwise conditional (in)dependencies among the variables in the graph, and the resulting model is therefore also called **conditional independence graph**.

If x_i and x_j two selected random variables/nodes, and the set $\{x_k\}$ represents all other variables/nodes with $k \neq i$ and $k \neq j$. We say that variables x_i and x_j are conditionally independent given all the other variables $\{x_k\}$

$$x_i \perp\!\!\!\perp x_j | \{x_k\}$$

if the joint probability density of x_i, x_j and x_k factorises as

$$\Pr(x_1, x_2, \dots, x_d) = \Pr(x_i | \{x_k\}) \Pr(x_j | \{x_k\}) \Pr(\{x_k\}).$$

or equivalently

$$\Pr(x_i, x_j | \{x_k\}) = \Pr(x_i | \{x_k\}) \Pr(x_j | \{x_k\}).$$

In a corresponding conditional independence graph, there is no edge between x_i and x_j , as in such a graph *missing edges correspond to conditional independencies* between the respective non-connected nodes.

5.2.5.1 Gaussian graphical model

Assuming that x_1, \dots, x_d are jointly normal distributed, i.e. $x \sim N(\mu, \Sigma)$, it turns out that it is straightforward to identify the pairwise conditional independencies. From Σ we first obtain the precision matrix

$$\Omega = (\omega_{ij}) = \Sigma^{-1}.$$

Crucially, it can be shown that $\omega_{ij} = 0$ implies $x_i \perp\!\!\!\perp x_j | \{x_k\}$! Hence, from the precision matrix Ω we can directly read off all the pairwise conditional independencies among the variables x_1, x_2, \dots, x_d !

Often, the covariance matrix Σ is dense (few zeros) but the corresponding precision matrix Ω is sparse (many zeros).

The conditional independence graph computed for normally distributed variables is called a **Gaussian graphical model**, or **GGM**. A further alternative name is **covariance selection model**.

5.2.5.2 Related quantity: partial correlation

From the precision matrix Ω we can also compute the matrix of pairwise full conditional *partial correlations*:

$$\rho_{ij|\text{rest}} = -\frac{\omega_{ij}}{\sqrt{\omega_{ii}\omega_{jj}}}$$

which is essentially the standardised precision matrix (similar to correlation but with an extra minus sign!)

The partial correlations lie in the range between -1 and +1, $\rho_{ij|\text{rest}} \in [-1, 1]$, just like standard correlations.

If x is multivariate normal then $\rho_{ij|\text{rest}} = 0$ indicates conditional independence between x_i and x_j .

Regression interpretation: partial correlation is the correlation that remains between the two variables if the effect of the other variables is “regressed away”. In other words, the partial correlation is exactly equivalent to the correlation between the residuals that remain after regressing x_i on the variables $\{x_k\}$ and x_j on $\{x_k\}$.

5.2.6 Algorithm for learning GGMs

From the above we can devise a simple algorithm to learn Gaussian graphical model (GGM) from data:

1. Estimate covariance $\hat{\Sigma}$ (in such a way that it is invertible!)
2. Compute corresponding partial correlations
3. If $\hat{\rho}_{ij|\text{rest}} \approx 0$ then there is (approx). conditional independence between x_i and x_j .
In practise this is done by statistical testing for vanishing partial correlations. If there are many edges we also need adjustment for simultaneous multiple testing since all edges are tested in parallel.

5.2.7 Example: exam score data (Mardia et al 1979:)

Correlations (rounded to 2 digits):

##	mechanics	vectors	algebra	analysis	statistics
## mechanics	1.00	0.55	0.55	0.41	0.39

```
## vectors      0.55    1.00    0.61    0.49    0.44
## algebra      0.55    0.61    1.00    0.71    0.66
## analysis     0.41    0.49    0.71    1.00    0.61
## statistics   0.39    0.44    0.66    0.61    1.00
```

Partial correlations (rounded to 2 digits):

```
##           mechanics vectors algebra analysis statistics
## mechanics    1.00    0.33    0.23    0.00    0.02
## vectors      0.33    1.00    0.28    0.08    0.02
## algebra      0.23    0.28    1.00    0.43    0.36
## analysis     0.00    0.08    0.43    1.00    0.25
## statistics   0.02    0.02    0.36    0.25    1.00
```

Note that there are no zero correlations but there are **four partial correlations close to 0**, indicating **conditional independencies** between:

- analysis and mechanics,
- statistics and mechanics,
- analysis and vectors, and
- statistics and vectors.

Thus, of 10 possible edges four are missing, and thus the conditional independence graph looks as follows:

```
Mechanics      Analysis
 |             \   /   |
 |      Algebra |
 |             /   \   |
 Vectors      Statistics
```


Chapter 6

Nonlinear and nonparametric models

In the last part of the module we discuss methods that go beyond the traditional linear methods in multivariate statistics.

Relevant textbooks:

The lectures for much of this part of the module follow selected chapters from the following three text books:

- James et al. (2013) *An introduction to statistical learning with applications in R*. Springer.
- Hastie et al. (2009) *The elements of statistical learning: data mining, inference, and prediction*. Springer.
- Rogers and Girolami (2017) *A first course in machine learning (2nd edition)*. CRC Press.

Please study the relevant section and chapters as indicated below in each subsection!

6.1 Limits of linear models and correlation

Linear models are very effective tools. However, it is important to recognise their limits especially when modelling complex nonlinear relationships.

6.1.1 Correlation only measures linear dependence

A very simple demonstration of this is given by the following example. Assume x is a normal distributed random variable with $x \sim N(0, 1)$. From x we construct a second random variable $y = x^2$ — thus y fully depends on x with no added extra noise. What is the correlation between x and y ?

Let's answer this question by doing a small computer simulation:

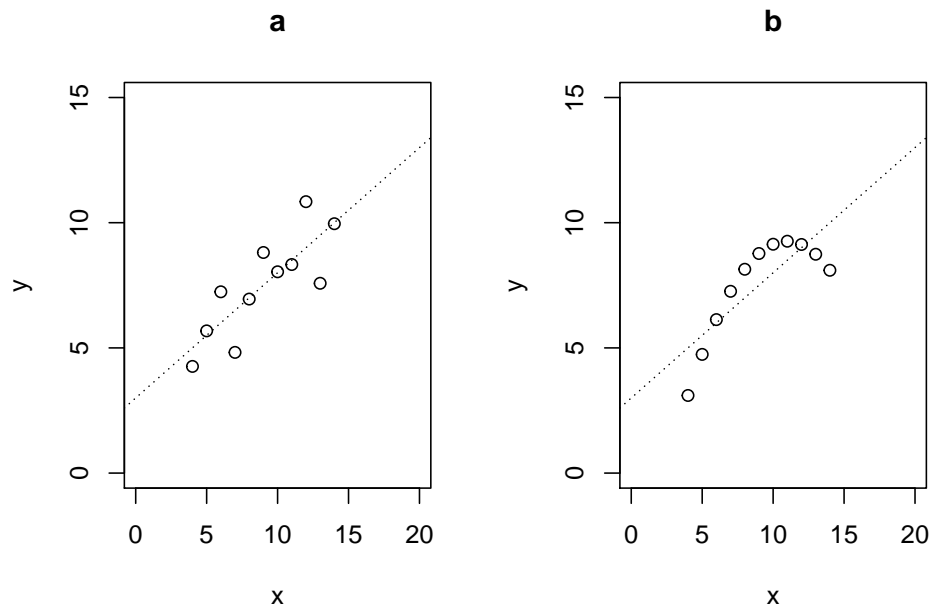
```
x=rnorm(10000)
y = x^2
cor(x,y)
```

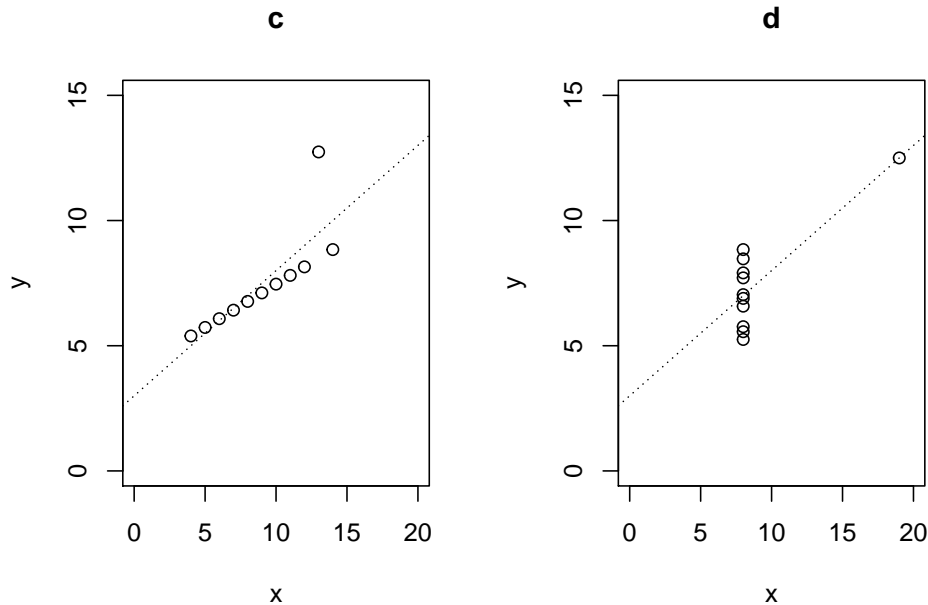
```
## [1] -0.01298635
```

Thus, correlation is (almost) zero even though x and y are full dependent! This is because correlation only measures linear dependence!

6.1.2 Anscombe data sets

Using correlation, and more generally linear models, blindly can thus hide complexities of the analysed data. A further classic example for this is demonstrated by the “Anscombe quartet” of data sets (F. J. Anscombe. 1973. Graphs in statistical analysis. The American Statistician 27:17-21, <http://dx.doi.org/10.1080/00031305.1973.10478966>):





As evident from the scatter plots the relationship between the two variables x and y is very different in the four cases! Intriguingly, all four data sets share exactly the same linear characteristics and summary statistics:

- Means $m_x = 9$ and $m_y = 7.5$
- Variances $s_x^2 = 11$ and $s_y^2 = 4.13$
- Correlation $r = 0.8162$
- Linear model fit with intercept $a = 3.0$ and slope $b = 0.5$

Thus, in actual data analysis it is always a **good idea to inspect the data visually** to get a first impression whether using a linear model makes sense.

In the above only data “a” follows a linear model. Data “b” represents a quadratic relationship. Data “c” is linear but with an outlier that disturbs the linear relationship. Finally data “d” also contains an outlier but also represent a case where y is (apart from the outlier) is not dependent on x .

In the Worksheet~10 a more recent version of the Anscomber quartet will be analysed in the form of the “datasauRus” dozen - 13 highly nonlinear datasets that all share the same linear characteristics.

6.2 Mutual information as generalised correlation

6.2.1 Definition of mutual information

Recall from the year 2 module “Statistical Methods” (MATH20802) the definition of Kullback-Leibler divergence, or relative entropy, between two distributions:

$$I^{KL}(F||G) := E_F \log \left(\frac{f(x)}{g(x)} \right)$$

Here F plays the role of the true distribution and G is an approximating distribution, with f and g being the corresponding densities.

The *Mutual information* (MI) between two random variables x and y is defined as the KL divergence between the corresponding joint distribution and the product distribution:

$$\text{MI}(x, y) = KL(F_{x,y} || F_x F_y) = E_{F_{x,y}} \log \left(\frac{f(x, y)}{f(x) f(y)} \right).$$

Thus, MI measures how well the joint distribution can be approximated by the product distribution (which would be the appropriate joint distribution if x and y are independent). Since MI is an application of KL divergence it shares all its properties. In particular, $\text{MI}(x, y) = 0$ implies that the joint distribution and product distributions are the same. Hence the two random variables x and y are independent if the mutual information vanishes.

6.2.2 Mutual information between two normal variables

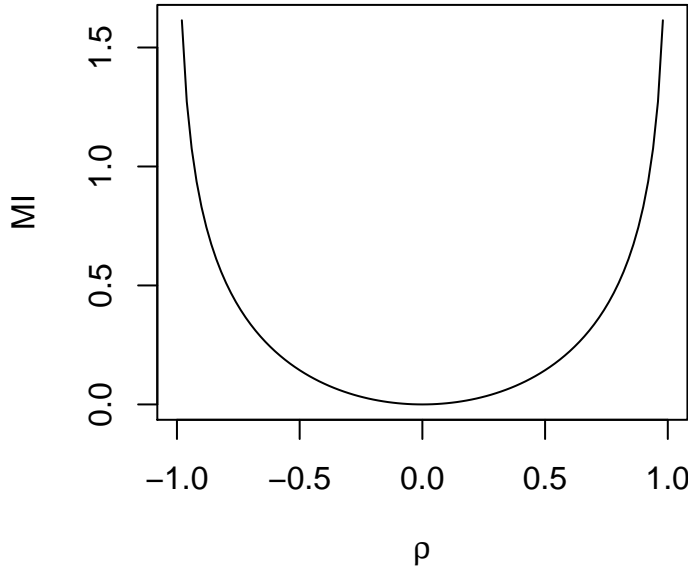
The KL divergence between two multivariate normal distributions F_0 and F is

$$I^{KL}(F_0 || F) = \frac{1}{2} \left\{ (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_0) + \text{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_0) - \log \det(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}_0) - d \right\}$$

This allows compute the mutual information $\text{MI}_{\text{norm}(x,y)}$ between two univariate random variables x and y that are correlated and assumed to be jointly bivariate normal. Let $\mathbf{z} = (x, y)^T$. The joint bivariate normal distribution is characterised by the mean $E(\mathbf{z}) = \boldsymbol{\mu} = (\mu_x, \mu_y)^T$ and the covariance matrix $\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}$ where $\text{Cor}(x, y) = \rho$. If x and y are independent then $\rho = 0$ and $\boldsymbol{\Sigma}_{\text{indep}} = \begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}$. The mutual information between x and y is therefore

$$\begin{aligned} \text{MI}_{\text{norm}}(x, y) &= I^{KL}(N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) || N(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{\text{indep}})) \\ &= \frac{1}{2} \left\{ \text{Tr}(\boldsymbol{\Sigma}_{\text{indep}}^{-1} \boldsymbol{\Sigma}) - \log \det(\boldsymbol{\Sigma}_{\text{indep}}^{-1} \boldsymbol{\Sigma}) - 2 \right\} \\ &= \frac{1}{2} \left\{ \text{Tr}(\mathbf{P}_{x,y}) - \log \det(\mathbf{P}_{x,y}) - 2 \right\} \\ &= -\frac{1}{2} \log \det(\mathbf{P}_{x,y}) \\ &= -\frac{1}{2} \log \det \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \\ &= -\frac{1}{2} \log(1 - \rho^2) \\ &\approx \frac{\rho^2}{2} \end{aligned}$$

Thus $\text{MI}_{\text{norm}}(x, y)$ is a one-to-one function of the squared correlation ρ^2 between x and y :



For small values of correlation $2\text{MI}_{\text{norm}}(x, y) \approx \rho^2$.

6.2.3 Mutual information between two normally distributed random vectors

The mutual information $\text{MI}_{\text{norm}}(x, y)$ between two multivariate normal random vector x and y can be computed in a similar fashion as in the bivariate case.

Let $z = (x, y)^T$ with dimension $d = p + q$. The joint multivariate normal distribution is characterised by the mean $E(z) = \mu = (\mu_x^T, \mu_y^T)^T$ and the covariance matrix $\Sigma = \begin{pmatrix} \Sigma_{x,x} & \Sigma_{x,y} \\ \Sigma_{x,y}^T & \Sigma_{y,y} \end{pmatrix}$. If x and y are independent then $\Sigma_{x,y} = 0$ and $\Sigma_{\text{indep}} = \begin{pmatrix} \Sigma_{x,x} & 0 \\ 0 & \Sigma_{y,y} \end{pmatrix}$. Mutual information between x and y is then

$$\begin{aligned}
 \text{MI}_{\text{norm}}(x, y) &= I^{KL}(N(\mu, \Sigma) || N(\mu, \Sigma_{\text{indep}})) \\
 &= \frac{1}{2} \left\{ \text{Tr} \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - \log \det \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - d \right\} \\
 &= -\frac{1}{2} \log \det \left(P_{\text{indep}}^{-1} P \right) \\
 &= -\frac{1}{2} \log \left(\det \begin{pmatrix} P_x & P_{xy} \\ P_{xy}^T & P_y \end{pmatrix} / (\det P_x \det P_y) \right) \\
 &= -\frac{1}{2} \log \det \left(I_q - P_y^{-1} P_{yx} P_x^{-1} P_{xy} \right) \\
 &= -\frac{1}{2} \log \det \left(I_p - P_x^{-1} P_{xy} P_y^{-1} P_{yx} \right)
 \end{aligned}$$

The equality of

$$\det \left(I_q - P_y^{-1} P_{yx} P_x^{-1} P_{xy} \right) = \det \left(I_p - P_x^{-1} P_{xy} P_y^{-1} P_{yx} \right)$$

follows from the Weinstein-Aronszajn determinant identity $\det(\mathbf{I} + \mathbf{AB}) = \det(\mathbf{I} + \mathbf{BA})$.

By comparison with the squared Rozeboom vector correlation coefficient $\rho_{x,y}^2$ (cf. Chapter 5) we see that

$$\text{MI}_{\text{norm}}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \log(1 - \rho_{x,y}^2) \approx \frac{1}{2} \rho_{x,y}^2$$

Thus, in the multivariate case $\text{MI}_{\text{norm}}(\mathbf{x}, \mathbf{y})$ has exactly the same functional relationship with the vector correlation $\rho_{x,y}^2$ as the MI for two univariate variables and squared Pearson correlation.

From the various expressions for the vector correlations we can also get further identities for the corresponding mutual information between \mathbf{x} and \mathbf{y} :

$$\begin{aligned} \text{MI}_{\text{norm}}(\mathbf{x}, \mathbf{y}) &= -\frac{1}{2} \log \det(\mathbf{I}_q - \mathbf{K}^T \mathbf{K}) \\ &= -\frac{1}{2} \log \det(\mathbf{I}_p - \mathbf{K} \mathbf{K}^T) \\ &= -\frac{1}{2} \sum_{i=1}^m \log(1 - \lambda_i^2) \end{aligned}$$

Note that the last line shows that $\text{MI}_{\text{norm}}(\mathbf{x}, \mathbf{y})$ is the sum of the MIs resulting from the individual canonical correlations (again with the same functional form as in the univariate Pearson correlation case).

6.2.4 Using MI for variable selection

In principle, MI can be computed for any distribution and model and thus applies to both normal and non-normal models, and to both linear and nonlinear relationships.

In very general way we may denote by $F_{y|x}$ we denote a predictive model for \mathbf{y} conditioned on \mathbf{x} and F_y is the marginal distribution of \mathbf{y} without predictors. Note that the predictive model can assume any form (incl. nonlinear). Typically $F_{y|x}$ is a complex model and F_y a simple model (no predictors).

Then mutual information between \mathbf{x} and \mathbf{y} can be also understood as expected KL divergence between the conditional and marginal distributions:

$$\mathbb{E}_{F_x} \text{KL}(F_{y|x} || F_y) = \text{MI}(\mathbf{x}, \mathbf{y})$$

This can be shown as follows. The KL-divergence between $F_{y|x}$ and F_y is given by

$$I^{KL}(F_{y|x}, F_y) = \mathbb{E}_{F_{y|x}} \log \left(\frac{f(\mathbf{y}|\mathbf{x})}{f(\mathbf{y})} \right),$$

which is a random variable since it depends on \mathbf{x} . Taking the expectation with regard to F_x (the distribution of \mathbf{x}) we get

$$\mathbb{E}_{F_x} I^{KL}(F_{y|x}, F_y) = \mathbb{E}_{F_x} \mathbb{E}_{F_{y|x}} \log \left(\frac{f(\mathbf{y}|\mathbf{x})f(\mathbf{x})}{f(\mathbf{y})f(\mathbf{x})} \right) = \mathbb{E}_{F_{x,y}} \log \left(\frac{f(\mathbf{x}, \mathbf{y})}{f(\mathbf{y})f(\mathbf{x})} \right) = \text{MI}(\mathbf{x}, \mathbf{y}).$$

Because of this link of MI with conditioning the MI between response and predictor variables is often used for variable and feature selection in general models.

6.2.5 Other measures of general dependence

Besides mutual information there are others measures of general dependence between multivariate random variables.

The two most important ones that have been proposed in the recent literature are i) distance correlation and ii) the maximal information coefficient (MIC and MIC_e).

6.3 Nonlinear regression models

Traditional linear (and generalised linear) models can be extended to nonlinear settings.

Relevant reading:

Please read: [James et al. \(2013\) Chapter 7 “Moving Beyond Linearity”](#)

Specifically:

- Section 7.1 Polynomial Regression
- Section 7.4 Regression Splines

6.3.1 Scatterplot smoothing

- lowess / loess algorithm

Locally weighted scatterplot smoothing (intended for exploratory analysis, not for probabilistic modelling).

6.3.2 Polynomial regression model

Advantage: - possible to use standard OLS tools to fit model and to do inference (relabeling trick for univariate models)

Disadvantage: - multivariate version complicated and intractable - high-order polynomials are very erratic - prone to overfitting if degree/order is too high

6.3.3 Piecewise polynomial regression

- simple linear piece-wise model
- basis function approach
- regression splines
- natural splines

See Worksheet~10 for practical application in R!

6.4 Random forests

Another widely used approach for prediction in nonlinear settings is the method of random forests.

Relevant reading:

Please read: [James et al. \(2013\)](#) Chapter 8 “Tree-Based Methods”

Specifically:

- Section 8.1 The Basics of Decision Trees
- Section 8.2.1 Bagging
- Section 8.2.2 Random Forests

6.4.1 Stochastic vs. algorithmic models

Two cultures in statistical modelling: stochastic vs. algorithmic models

Classic discussion paper by Leo Breiman (2001): Statistical modeling: the two cultures. Statistical Science. Vol 16, pages 199-231. <https://projecteuclid.org/euclid.ss/1009213726>

6.4.2 Random forests

Invented by Breimann in 1996.

Basic idea:

- A single decision tree is unreliable and unstable (weak predictor/classifier).
- Use bootstrap to generate multiple decision trees (=“forest”)
- Average over predictions from all tree (=“bagging”, bootstrap aggregation)

The averaging procedure has the effect of variance stabilisation. Intriguingly, averaging across all decision trees dramatically improves the overall prediction accuracy!

The Random Forests approach is an example of an **ensemble method** (since it is based on using an “ensemble” of trees).

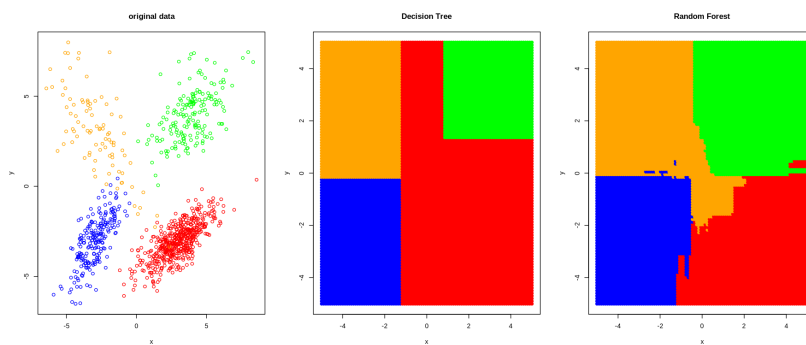
Variations: boosting, XGBoost (<https://xgboost.ai/>)

Random forests will be applied in Computer Lab 5.

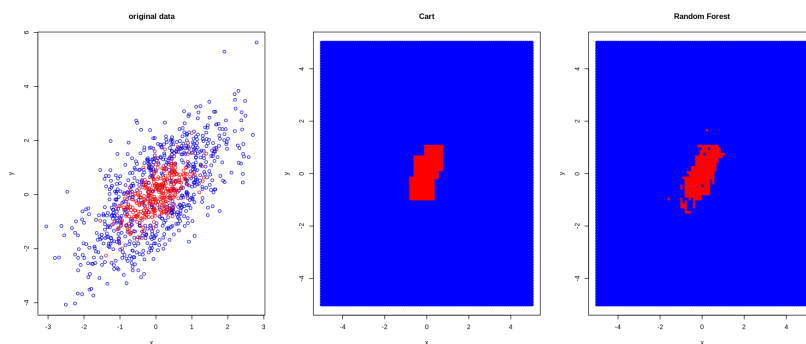
They are computationally expensive but typically perform very well!

6.4.3 Comparison of decision boundaries: decision tree vs. random forest

Non-nested case:



Nested case:



Compare also with the decision boundaries for LDA and QDA (previous chapter).

See Worksheet~11 for practical application of random forests in R!

6.5 Gaussian processes

Gaussian processes offer another nonparametric approach to model nonlinear dependencies. They provide a probabilistic model for the unknown nonlinear function.

Relevant reading

Please read: [Rogers and Girolami \(2017\) Chapter 8: Gaussian processes.](#)

6.5.1 Main concepts

- Gaussian processes (GPs) belong to the family of **Bayesian nonparametric models**
- Idea:
 - start with prior over a function (!),
 - then condition on observed data to get posterior distribution (again over all functions)
 - use an infinitely dimensional multivariate normal distribution as prior

6.5.2 Technical background:

GPs make use of the fact that marginal and conditional distributions of a multivariate normal are also multivariate normal.

Multivariate normal distribution:

$$z \sim N_d(\mu, \Sigma)$$

Assume:

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

with

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

and

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{pmatrix}$$

with corresponding dimensions d_1 and d_2 and $d_1 + d_2 = d$.

Marginal distributions:

Any subset of z is also multivariate normal distributed:

$$z_i \sim N_{d_i}(\mu_i, \Sigma_{ii})$$

Conditional multivariate normal:

The conditional distribution is also multivariate normal:

$$z_i | z_j = z_{i|j} \sim N_{d_i}(\mu_{i|j}, \Sigma_{i|j})$$

with

$$\mu_{i|j} = \mu_i + \Sigma_{ij} \Sigma_{jj}^{-1} (z_j - \mu_j)$$

and

$$\Sigma_{i|j} = \Sigma_{ii} - \Sigma_{ij} \Sigma_{jj}^{-1} \Sigma_{ij}^T$$

$z_{i|j}$ and $\mu_{i|j}$ have dimension $d_i \times 1$ and $\Sigma_{i|j}$ has dimension $d_i \times d_i$

6.5.3 Covariance functions and kernel

The GP prior is a infinitely dimensional multivariate normal with mean zero and the **covariance specified by a function**:

A widely used covariance function is

$$\text{Cov}(x, x') = \sigma^2 e^{-\frac{(x-x')^2}{2l^2}}$$

This is known as the **squared-exponential kernel** or **Radial-basis function (RBF) kernel**.

Note that $(x, x) = \sigma^2$ and the autocorrelation $\text{Cor}(x, x') = e^{-\frac{(x-x')^2}{2l^2}}$.

The parameter l is the length scale parameter and describes the wigglyness or smoothness of the resulting function. Small values of l mean more complex, more wiggly functions, and low autocorrelation.

There are many other kernel functions, including periodic, polynomial and linear kernels.

6.5.4 GP model

Nonlinear regression in the GP approach is conceptually very simple:

- start with GP prior over all x
- then condition on the observed x_1, \dots, x_n
- the resulting conditional multivariate normal can be used to predict the function values at any unobserved values of x
- automatically provides credible intervals for predictions.

GP regression also provides a direct link with Bayesian linear regression (using a linear kernel).

Drawbacks: computationally expensive (n^3 because of the matrix inversion)

6.6 Neural networks

Another highly important class of models for nonlinear prediction (and nonlinear function approximation) are neural networks.

Relevant reading:

Please read: [Hastie et al. \(2009\) Chapter 11 “Neural networks”](#)

6.6.1 History

Neural networks are actually relatively old models, going back to the 1950s!

Three phases of neural networks (NN)

- 1950/60: replicating functions of neurons in the brain (perceptron)
- 1980/90: neural networks as universal function approximators
- 2010—today: deep learning

The first phase was biologically inspired, the second phase focused on mathematical properties, and the current phase is pushed forward by advances in computer science and numerical optimisation:

- backpropagation algorithm
- auto-differentiation,
- stochastic gradient descent
- use of GPUs and TPUs,
- availability and development of software packages by major internet companies:
 - TensorFlow/Keras (Google),

- MXNet (Amazon),
- PyTorch (Facebook),
- PaddlePaddle (Baidu) etc.

6.6.2 Neural networks

Neural networks are essentially stacked systems of linear regressions, mapping input nodes (random variables) to outputs (response nodes). Each internal layer corresponds to internal latent variables. Each layer is connected with the next layer by **non-linear activation functions**.

- feedforward single layer NN
- stacked nonlinear multiple regression with hidden variables
- optimise by empirical risk minimisation

It can be shown that NN can approximate any arbitrary non-linear function mapping input and output.

“Deep” neural networks have many layers, and their optimisation requires advanced techniques (see above).

Neural networks are very highly parameterised models and require typically a lot of data for training.

Some of the statistical aspects of NN are not well understood: in particular it is known that NN overfit the data but can still generalise well. On the other hand, it is also known that NN can also be “fooled”, i.e. prediction can be unstable (adversarial examples).

Current statistical research on NN focuses on interpretability and on links with Bayesian inference and models (e.g. GPs). For example:

- <https://link.springer.com/book/10.1007/978-3-030-28954-6>
- <https://arxiv.org/abs/1910.12478>

6.6.3 Learning more about deep learning

A good place to learn more about deep learning and about the actual implementations in computer code on various platforms is the book “Dive into deep learning” by Zhang et al. (2020) available online at <https://d2l.ai/>

Appendix A

Brief refresher on matrices

This is intended a very short recap of some essentials you need to know about matrices. For more details please consult the lecture notes of earlier modules (e.g. linear algebra).

A.1 Matrix notation

We will frequently make use of matrix calculations. Matrix notation helps to make the equations simpler and enables to understand them better.

In matrix notation we distinguish between scalars, vectors, and matrices:

Scalar: x , X , lower or upper case, plain type.

Vector: \mathbf{x} , lower case, bold type. In handwriting one uses an arrow \vec{x} to indicate a vector.

Matrix: \mathbf{X} , upper case, bold type. In handwriting you use an underscore \underline{X} to indicate a matrix.

Note that a vector may be viewed as a matrix (with only one column or only one row). Likewise, a scalar may also be considered as a special case of a matrix.

Note on **random** matrices and vectors:

In the above notation you need to determine from the context whether a quantity represents a random variable, or whether it is a constant. You cannot read this from the case (upper vs. lower case) as in the standard notation commonly used in univariate statistics.

A.2 Simple special matrices

I_d is the identity matrix. It is a square matrix of dimension $d \times d$ with the diagonal filled with 1 and off-diagonals filled with 0.

$$I_d = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & & 0 \\ \vdots & \vdots & & \ddots & \\ 0 & 0 & 0 & & 1 \end{pmatrix}$$

$\mathbf{1}$ is a matrix that contains only 1s. Most often it is used in the form of a column vector with d rows:

$$\mathbf{1}_d = \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

A.3 Simple matrix operations

Matrices behave much like ordinary numbers. For example, you can apply operations such as matrix addition $A + B$ and matrix multiplication AB . In order to conduct these operations the matrices need to be compatible: for addition the matrices need to have the same dimension, and for multiplication the number of columns of A must match the number of rows of B . Note that $AB \neq BA$, i.e. matrix multiplication is in general not commutative.

If A is a squared matrix and is nonsingular (i.e. it has no zero eigenvalues) then you can define an inverse A^{-1} such that $A^{-1}A = AA^{-1} = I$.

The matrix transpose $t(A) = A^T$ interchanges rows and columns.

The trace of the matrix is the sum of the diagonal entries $\text{Tr}(A) = \sum a_{ii}$.

The sum of the squares of all entries of a rectangular matrix $A = (a_{ij})$ can be written using the trace as follows: $\sum_{i,j} a_{ij}^2 = \text{Tr}(A^T A) = \text{Tr}(AA^T)$

A.4 Orthogonal matrices

An orthogonal matrix Q has the property that $Q^T = Q^{-1}$. This implies that $QQ^T = Q^T Q = I$. An orthogonal matrix Q can be interpreted geometrically as a rotation-reflection operator since multiplication of Q with a vector will result in a new vector of the same length but with a change in direction. The identity matrix I is the simplest example of an orthogonal matrix but in general there are infinitely many rotation-reflection matrices.

A.5 Eigenvalues and eigenvalue decomposition

A vector \mathbf{u}_i is called an eigenvector of a square matrix A and λ_i the corresponding eigenvalue if $A\mathbf{u}_i = \mathbf{u}_i\lambda_i$.

If A is a symmetric matrix (i.e. $A = A^T$) with real entries (as assumed for the rest of this section) then it can be shown that all eigenvalues are real, and that the corresponding eigenvectors are all orthogonal.

The eigenvalue decomposition of a symmetric real-valued A is given by

$$A = U\Lambda U^T$$

with

$$\Lambda = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_d \end{pmatrix}$$

being a diagonal matrix containing all eigenvalues λ_i of A . U is an orthogonal matrix containing the corresponding eigensystem (i.e. all eigenvectors \mathbf{u}_i) in the columns of U . The eigenvectors in U are orthonormal, i.e. they are orthogonal to each other and have length 1.

If A is multiplied with U we immediately see that Λ contains the eigenvalues since $AU = U\Lambda$.

Furthermore, it can be shown that the above eigenvalue decomposition of A is **unique apart from the signs of the eigenvectors** (i.e. individual column signs of U can be changed and the eigenvalue decomposition is still valid). In order to make the eigenvalue decomposition fully unique you need to impose further restrictions (e.g. require a positive diagonal of U). Note that this may be particularly important in computer applications where the sign can vary depending on the specific implementation of the underlying numerical algorithms.

Eigenvalue decomposition is also known as spectral decomposition.

A.6 Singular value decomposition

A generalisation of the eigenvalue decomposition (which is for squared matrices) is the **singular value decomposition** (SVD).

Let A be a $n \times m$ matrix. The SVD decomposition of A is $A = UDV^T$.

U and V are orthogonal matrices, D contains the singular values.

As the eigenvalue decomposition SVD is unique apart from the signs of the column vectors in U and V .

A.7 Positive (semi-)definiteness, rank, condition

If all $\lambda_i \geq 0$ then A is called positive semi-definite.

If all eigenvalues are strictly positive $\lambda_i > 0$ then A is called positive definite.

If one or more of the eigenvalues equal zero then A is said to be singular.

The rank of A counts how many eigenvalues are non-zero.

A is full rank if all eigenvalues are non-zero.

The condition number of A is the ratio between the largest and smallest absolute eigenvalue.

If A is singular then the condition number is infinite.

A.8 Trace and determinant of a matrix and eigenvalues

The trace of the matrix A can be also obtained as the *sum* of its eigenvalues: $\text{Tr}(A) = \sum \lambda_i$.

The determinant of A is the *product* of the eigenvalues, $\det(A) = \prod \lambda_i$. Therefore, if A is singular then $\det(A) = 0$.

Determinants have a multiplicative property, $\det(AB) = \det(A) \det(B)$. Another important identity is $\det(I_n + AB) = \det(I_m + BA)$ where A is a $n \times m$ and B is a $m \times n$ matrix. This is called the Weinstein-Aronszajn determinant identity (also credited to Sylvester).

A.9 Functions of matrices

Again, we focus on symmetric, real-valued squared matrices A .

A matrix function $f(A)$, generalising from a simple function $f(a)$ can then be defined via the eigenvalue decomposition as

$$f(A) = U f(\Lambda) U^T = U \begin{pmatrix} f(\lambda_1) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & f(\lambda_d) \end{pmatrix} U^T$$

Therefore, in order to obtain the corresponding matrix function, the function is simply applied on the level of eigenvalues. By construction $f(A)$ will also be symmetric and real-valued as long as the transformed eigenvalues $f(\lambda_i)$ are real.

Examples:

Example A.1. matrix power: $f(a) = a^p$ (with p real number)

Special cases of matrix power include :

- matrix inversion: $f(a) = a^{-1}$
- the matrix square root: $f(a) = a^{1/2}$
(since there are multiple solutions to the square root there are also multiple matrix square roots. The principal matrix square root is given by using the positive square roots of all the eigenvalues. Thus the **principal matrix square root** of a positive semidefinite matrix is also positive semidefinite and it is unique).

Example A.2. matrix exponential: $f(a) = \exp(a)$

Example A.3. matrix logarithm: $f(a) = \log(a)$

For a positive definite matrix A computing the trace of the matrix logarithm of A is the same as taking the log of the determinant of A :

$$\text{Tr}(\log(A)) = \log \det(A)$$

because $\sum \log(\lambda_i) = \log(\prod \lambda_i)$.

A.10 Matrix calculus

A.10.1 First order vector derivatives

A.10.1.1 Gradient

The **nabla operator** (also known as **del operator**) is the *row* vector

$$\nabla = \left(\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_d} \right) = \frac{\partial}{\partial \mathbf{x}}$$

containing the first order partial derivative operators.

The **gradient** of a scalar-valued function $f(\mathbf{x})$ with vector argument $\mathbf{x} = (x_1, \dots, x_d)^T$ is also a *row* vector (with d columns) and can be expressed using the nabla operator

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_d} \right) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \text{grad} f(\mathbf{x}).$$

Note the various notations for the gradient.

Example A.4. $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$. Then $\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}^T$.

Example A.5. $f(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$. Then $\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{x}^T$.

Example A.6. $f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$. Then $\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{x}^T (A + A^T)$.

A.10.1.2 Jacobian matrix

For a vector-valued function

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^T.$$

the computation of the gradient of each component yields the **Jacobian matrix** (with m rows and d columns)

$$J_f(\mathbf{x}) = \begin{pmatrix} \nabla f_1(\mathbf{x}) \\ \vdots \\ \nabla f_m(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_j} \\ \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_j} \end{pmatrix} = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = D\mathbf{f}(\mathbf{x})$$

Again, note the various notations for the Jacobian matrix!

Example A.7. $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$. Then $J_f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = A$.

If $m = d$ then the Jacobian matrix is a square matrix and this allows to compute the **Jacobian determinant**

$$\det J_f(\mathbf{x}) = \det \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)$$

If $\mathbf{y} = f(\mathbf{x})$ is an invertible function with $\mathbf{x} = f^{-1}(\mathbf{y})$ then the Jacobian matrix is invertible and the inverted matrix is in fact the Jacobian of the inverse function!

This allows to compute the Jacobian determinant of the backtransformation as the inverse of the Jacobian determinant the original function:

$$\det Df^{-1}(\mathbf{y}) = (\det Df(\mathbf{x}))^{-1}$$

or in alternative notation

$$\det D\mathbf{x}(\mathbf{y}) = \frac{1}{\det D\mathbf{y}(\mathbf{x})}$$

A.10.2 Second order vector derivatives

The matrix of all second order partial derivatives of scalar-valued function with vector-valued argument is called the **Hessian matrix** and is computed by double application of the nabla operator:

$$\nabla^T \nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_d \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_d \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_d^2} \end{pmatrix} = \left(\frac{\partial f(\mathbf{x})}{\partial x_i \partial x_j} \right) = \left(\frac{\partial}{\partial \mathbf{x}} \right)^T \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}.$$

By construction it is square and symmetric.

A.10.3 First order matrix derivatives

The derivative of a scalar-valued function $f(\mathbf{X})$ with regard to a matrix argument \mathbf{X} can also be defined and results in a matrix with transposed dimensions compared to \mathbf{X} .

Two important specific examples are:

Example A.8. $\frac{\partial \text{Tr}(A\mathbf{X})}{\partial \mathbf{X}} = A$

Example A.9. $\frac{\partial \log \det(\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{Tr}(\log \mathbf{X})}{\partial \mathbf{X}} = \mathbf{X}^{-1}$

Appendix B

Further study

In this module we can only touch the surface of the field of multivariate statistics and machine learning. If you would like to study further I recommend the following books below as a starting point.

B.1 Recommended reading

For multivariate statistics and machine learning:

- Härdle and Simar (2015) *Applied multivariate statistical analysis. 4th edition.* Springer.
- Hastie et al. (2009) *The elements of statistical learning: data mining, inference, and prediction.* Springer.
- James et al. (2013) *An introduction to statistical learning with applications in R.* Springer.
- Marden (2015) *Multivariate Statistics: Old School*
- Rogers and Girolami (2017) *A first course in machine learning (2nd Edition).* Chapman and Hall / CRC.

B.2 Advanced reading

Additional (advanced) reference books for probabilistic machine learning are:

- Murphy (2012) *Machine learning: a probabilistic perspective.* MIT Press.
- Bishop (2006) *Pattern recognition and machine learning.* Springer.

Bibliography

- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer. <https://www.microsoft.com/en-us/research/people/cmbishop/prml-book/>.
- Härdle, W. K. and Simar, L. (2015). *Applied Multivariate Statistical Analysis*. Springer, Berlin.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2nd edition. <https://web.stanford.edu/~hastie/ElemStatLearn/>.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning with Applications in R*. Springer. <http://faculty.marshall.usc.edu/gareth-james/ISL/>.
- Marden, J. I. (2015). *Multivariate Statistics: Old School*. CreateSpace. <http://stat.istics.net/Multivariate>.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
- Rogers, S. and Girolami, M. (2017). *A first course in machine learning*. Chapman and Hall / CRC, 2nd edition.
- Zhang, A., Lipton, Z. C., Li, M., and Smola, A. J. (2020). *Dive into Deep Learning*. <https://d2l.ai>.